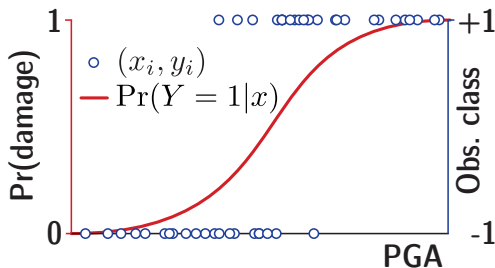


What is classification?



Data

$$\mathcal{D} = \{(x_i, y_i), \forall i = 1 : D\}$$

$$x_i \in \mathbb{R} : \begin{cases} \text{Covariate} \\ \text{attribute} \\ \text{regressor} \end{cases}$$

$$y_i \in \{-1, 1\} : \text{Observation}$$

Classification methods: mathematical models for $\Pr(Y|x)$

2 Types of classification methods

There are two types of classification methods:

1. **Generative**

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$














2. **Discriminative**

$$p(y|x) = g(x)$$

📍 Module #6 Outline

Introduction
Generative classifier
Logistic Regression
GPC
Neural Networks
Summary

Topics organization

<i>Background</i>	1	Revision probability & linear algebra	
	2	Probability distributions	 
<i>Machine Learning Basics</i>	0	Introduction	
	3	Bayesian Estimation	$\rho(A B) = \frac{\rho(B A)\rho(A)}{\rho(B)}$
	4	MCMC sampling & Newton	 
<i>Supervised learning</i>	5	Regression	
	6	Classification	 
	7	LSTM networks for time series	
<i>Unsupervised learning</i>	7	State-space model for time-series	
<i>Decision Making & RL</i>	8	Decision Theory	
	9	AI & Sequential decision problems	

Section Outline

Generative classifier

- 2.1 Introduction
 - 2.2 Formulation
 - 2.3 Multiple classes
 - 2.4 Example
 - 2.5 Classic generative methods
-

Formulation - Binary classification $y \in \{-1, 1\}$

Posterior probability of a class

$$p(y|x) = \begin{cases} \Pr(Y = -1|x) \\ \Pr(Y = +1|x) = 1 - \Pr(Y = -1|x) \end{cases}$$

Prior probability of a class

$$p(y) = \begin{cases} \Pr(Y = -1) \\ \Pr(Y = +1) \end{cases}$$

$$\underbrace{p(y|x)}_{\text{posterior}} = \frac{\underbrace{f(x|y)}_{\text{likelihood}} \cdot \underbrace{p(y)}_{\text{prior}}}{\underbrace{f(x)}_{\text{norm. cte.}}}$$

Likelihood of x_i given a class

$$f(x|y) = \begin{cases} f(x|y = -1) \quad \text{📈} \\ f(x|y = +1) \quad \text{📈} \end{cases}$$

Normalization constant

$$f(x) = \sum_{y \in \{-1, 1\}} f(x|y) \cdot p(y) = \text{cte.}$$

Learning parameters - Binary classification

$$\underbrace{p(y|x)}_{\text{posterior}} = \frac{\overbrace{f(x|y)}^{\text{likelihood}} \cdot \overbrace{p(y)}^{\text{prior}}}{\underbrace{f(x)}_{\text{norm. cte.}}}$$

**We need to learn the PDF/PMF parameters from \mathcal{D}
MLE/Bayes**

Formulation - Special case $f(x|y) = \mathcal{N}(x; \mu, \sigma^2)$

If we assume that $f(x|y=j) = \mathcal{N}(x; \mu_{x|y_j}^*, \sigma_{x|y_j}^{2*})$, and that the number of available data \mathbf{D} is large

MLE approximation of $\mu_{x|y_j}$ & $\sigma_{x|y_j}^2$

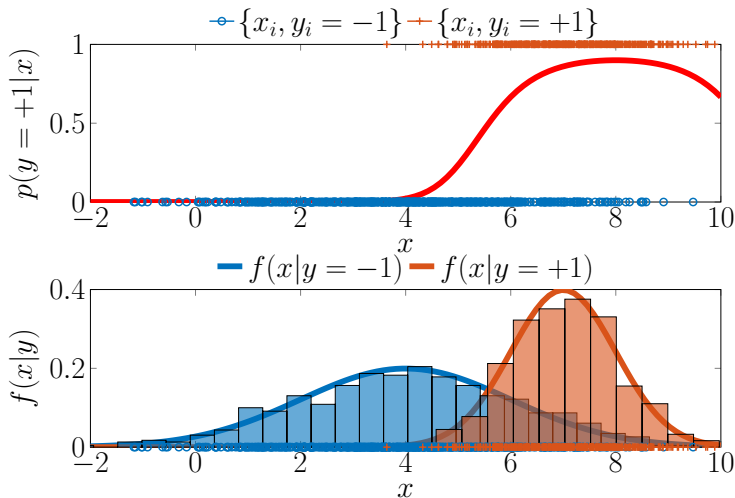
$$\mu_{x|y_j}^* = \frac{1}{\#\{i:y_i=j\}} \sum_i x_i, \quad \forall \{i : y_i = j\}$$

$$\sigma_{x|y_j}^{2*} = \frac{1}{\#\{i:y_i=j\}} \sum_i (x_i - \mu_{x|y_j}^*)^2, \quad \forall \{i : y_i = j\}$$

MLE approximation of $p(y)$

$$p^*(y) = \left\{ \Pr^*(y = j) = \frac{\#\{i:y_i=j\}}{\mathbf{D}} \right.$$

Example – Binary classification



[CIV_ML/Classification_example_1.m]

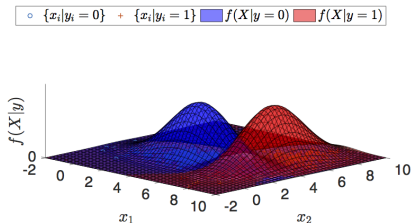
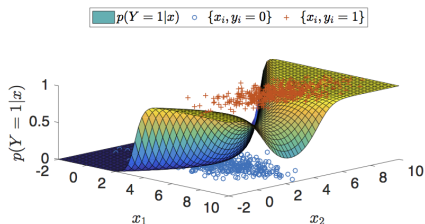
Multiattribute classification

What if there is **more than one attribute** $\mathbf{x} = [x_1, x_2, \dots, x_X]$?

Employ joint PDFs: $f(\mathbf{x}|y = j)$

$$f(\mathbf{x}|y = j) = \mathcal{N}(\mathbf{x}; \underbrace{\mu_{\mathbf{x}|y_j}, \Sigma_{\mathbf{x}|y_j}}_{\text{MLE/Bayes}})$$

Example – Multiattribute classification



Formulation - Multi-classes classification

$$y \in \{1, 2, \dots, Y\}$$

$$\underbrace{p(y|x)}_{\text{posterior}} = \frac{\underbrace{f(x|y)}_{\text{likelihood}} \cdot \underbrace{p(y)}_{\text{prior}}}{\underbrace{f(x)}_{\text{norm. cte.}}}$$

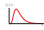
Posterior probability of a class

$$p(y|x) = \{ \Pr(Y = j|x) \}$$

Prior probability of a class

$$p(y) = \{ \Pr(Y = j) \}$$

Likelihood of x_i given a class

$$f(x|y) = \{ f(x|y = j) \}$$


Normalization constant

$$f(x) = \sum_j f(x|y = j) \cdot p(y = j)$$

Context

Societal challenge

Right after an earthquake

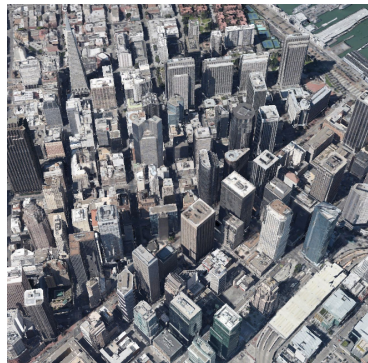
- which building is safe?
- which one should be evacuated?

Solution

Monitor structures

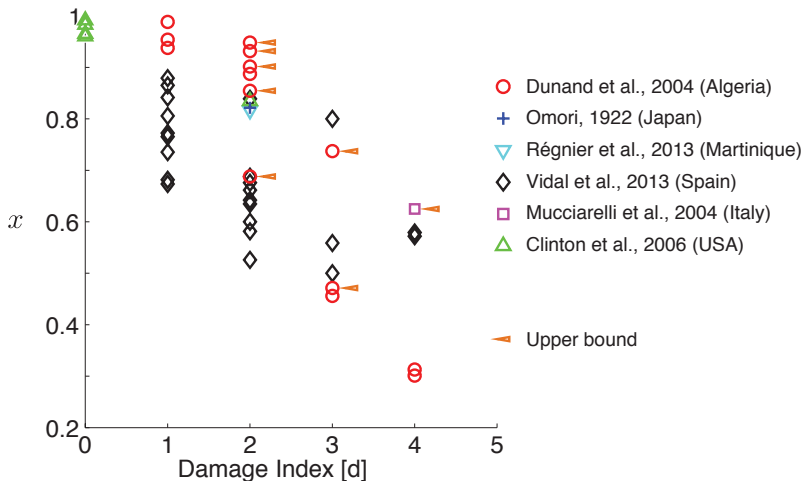
Scientific challenges

Real-time data-driven learning

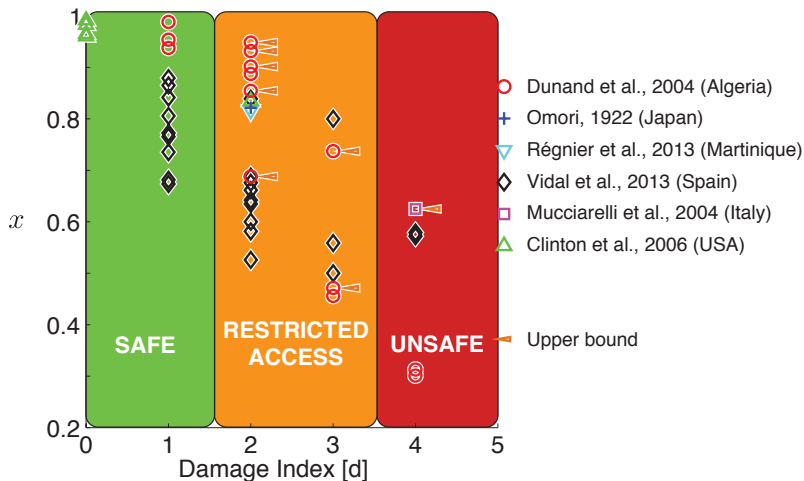


[Google images]

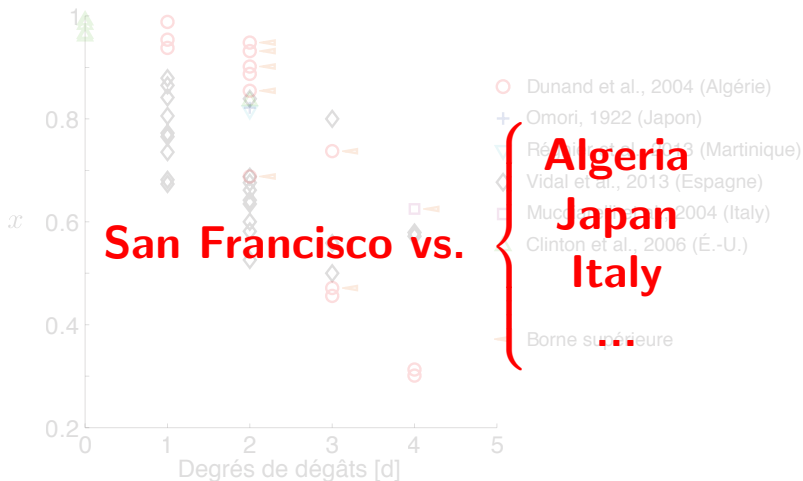
Frequency ratio: $x = f_{q,\text{post-earthquake}} / f_{q,\text{pre-earthquake}}$



Frequency ratio: $x = f_{q \cdot \text{post-earthquake}} / f_{q \cdot \text{pre-earthquake}}$



Frequency ratio: $x = f_{q \cdot \text{post-earthquake}} / f_{q \cdot \text{pre-earthquake}}$



Frequency ratio: $x = f_{q \cdot \text{post-earthquake}} / f_{q \cdot \text{pre-earthquake}}$

**We need to learn
as inspections are realized
after an earthquake**

Step 1 - Probabilistic learning

Given: $\mathcal{D} = \left\{ \underbrace{x_i}_{\in(0,1)}, \underbrace{d_i}_{\in\{0:5\}} \right\}, \forall i = 1 : D$

Goal: Learn what is the conditional PDF (predictive) of X given D

$$f(x; \theta(d)) = \mathcal{B}(x; \underbrace{\mu(d), \sigma(d)}_{\theta})$$

$$\underbrace{f(\theta(d)|\mathcal{D})}_{\text{posterior}} = \frac{\underbrace{f(\mathcal{D}|\theta(d))}_{\text{likelihood}} \cdot \underbrace{f(\theta(d))}_{\text{prior}}}{\underbrace{f(\mathcal{D})}_{\text{cte.}}}$$

$$f(x|d, \mathcal{D}) = \int f(x; \theta(d)) \cdot f(\theta(d)|\mathcal{D}) d\theta$$

Step 1 - Probabilistic learning (likelihood & prior)

$$\overbrace{f(\mu(d), \sigma(d) | \mathcal{D})}^{\text{posterior}} \propto \overbrace{f(\mathcal{D} | \mu(d), \sigma(d))}^{\text{likelihood}} \cdot \overbrace{f(\mu(d), \sigma(d))}^{\text{prior}}$$

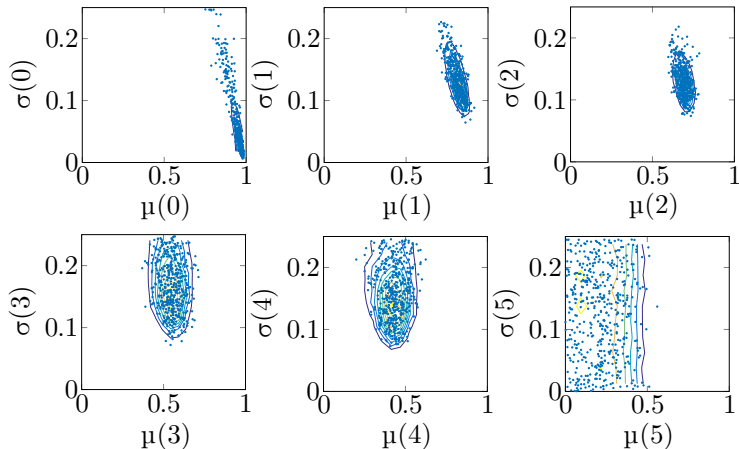
$$f(\mathcal{D} | \theta) = \prod_{\{i: d_i = d\}} \mathcal{L}(x_i | \mu(d), \sigma(d))$$

$$\mathcal{L}(x_i | \mu(d), \sigma(d)) = \begin{cases} f(x_i | \mu(d), \sigma(d)), & \text{if } x_i \text{ direct obs.} \\ F(x_i | \mu(d), \sigma(d)), & \text{if } x_i \text{ censored obs.} \end{cases}$$

We assume no prior information except for the constraint that

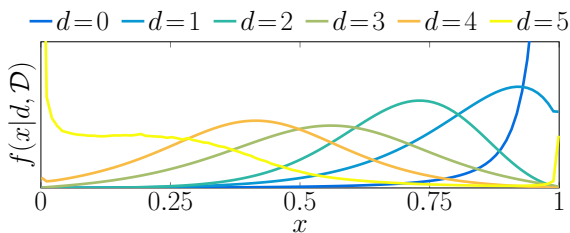
$$\mu(0) > \mu(1) > \dots > \mu(5)$$

Application - $f(\mu(d), \sigma(d)|\mathcal{D})$ & $f(x|d)$



Application - $f(\mu(d), \sigma(d)|\mathcal{D})$ & $f(x|d)$

$$f(x|d, \mathcal{D}) = \iint \mathcal{B}(x; \mu(d), \sigma(d)) f(\mu(d), \sigma(d)|\mathcal{D}) d\mu(d) d\sigma(d)$$

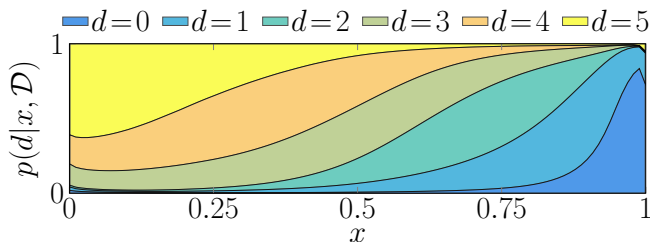


$$f(x|\theta(d)) = \mathcal{B}(x; \mu(d), \sigma(d))$$

We have $f(x|d, \mathcal{D})$, we want $p(d|x, \mathcal{D})$

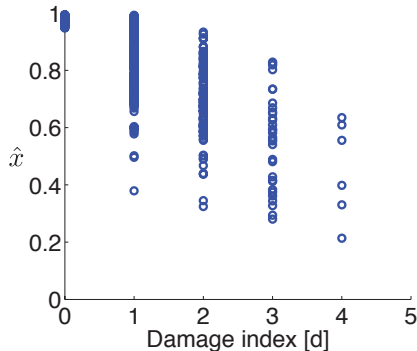
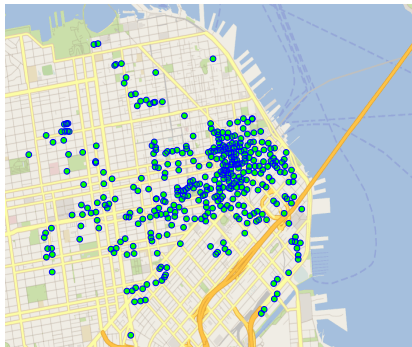
Step 2- Probabilistic prognosis

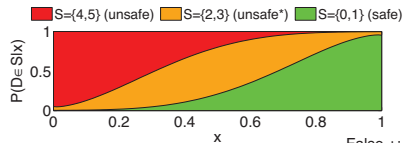
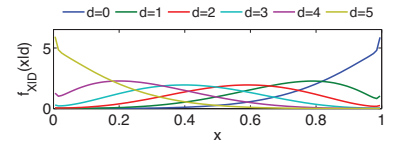
$$p(d|x, \mathcal{D}) = \int \frac{f(x; \theta(d)) \cdot p(d)}{\sum_{d'=0}^5 f(x; \theta(d')) \cdot p(d')} \cdot f(\theta(d)|\mathcal{D}) d\theta$$



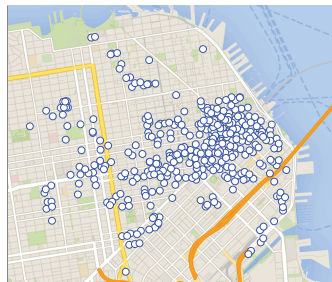
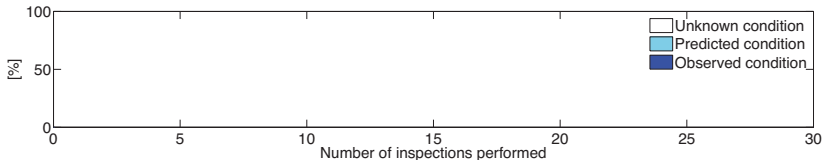
Post-earthquake damage simulation for San Francisco

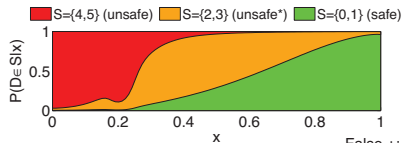
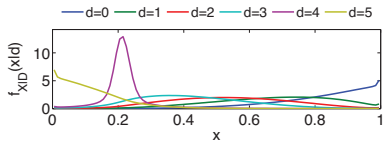
389 buildings ≥ 10 story



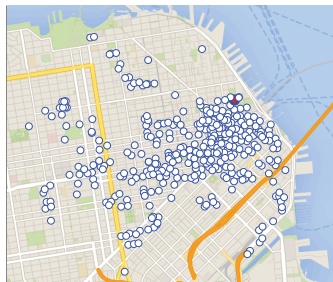
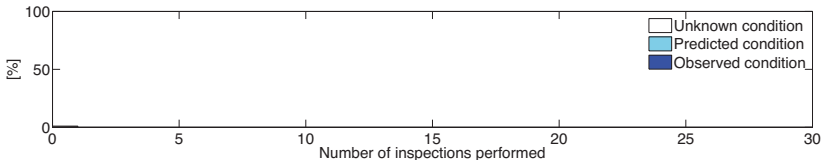


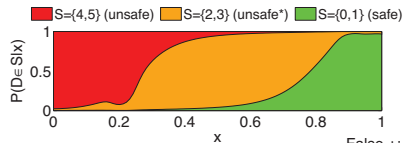
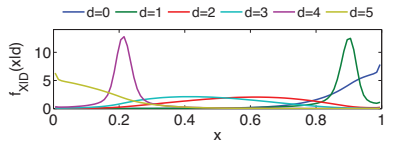
False +: 0% | False -: 0%



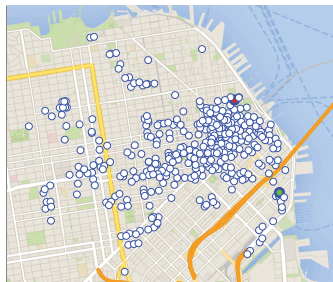
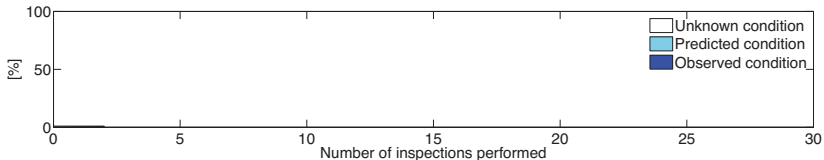


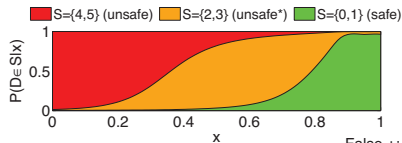
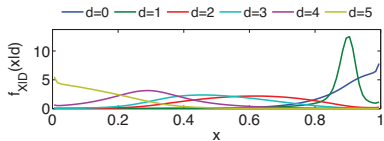
False +: 0% | False -: 0%



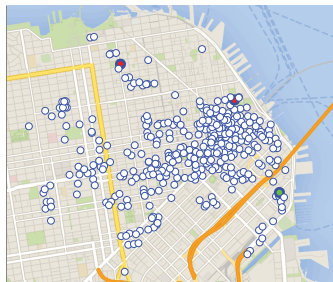
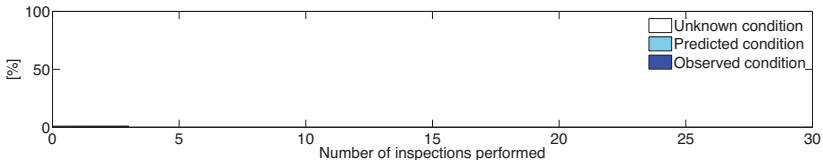


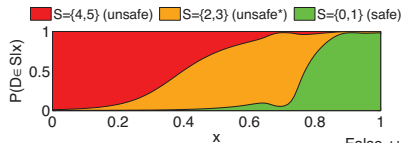
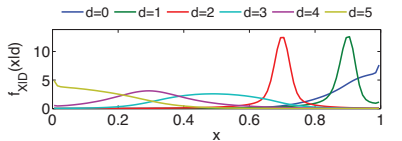
False +: 0% | False -: 0%



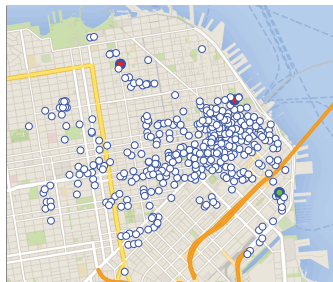


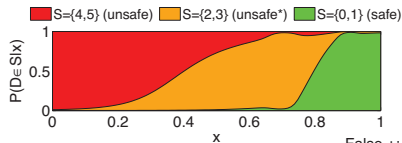
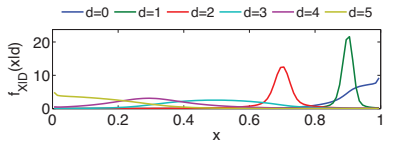
False +: 0% | False -: 0%



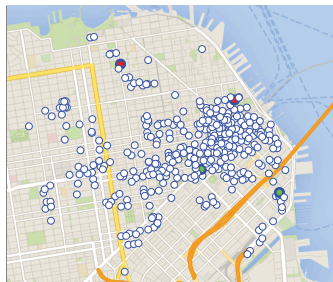


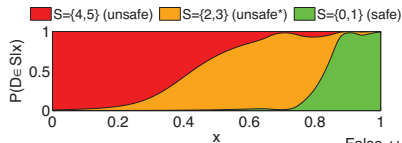
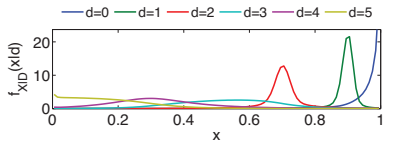
False +: 0% | False -: 0%



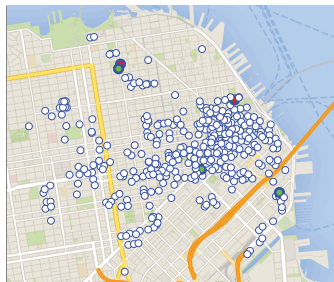
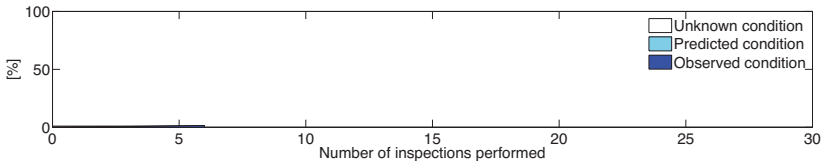


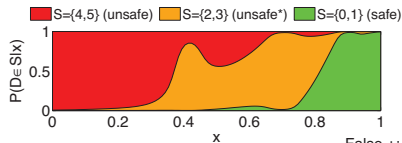
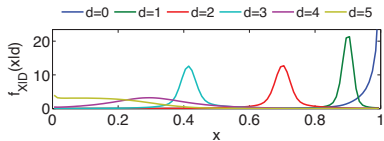
False +: 0% | False -: 0%



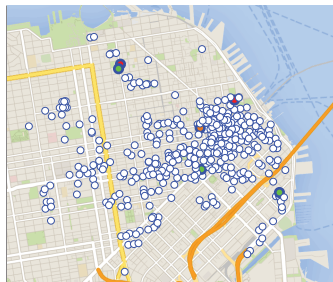
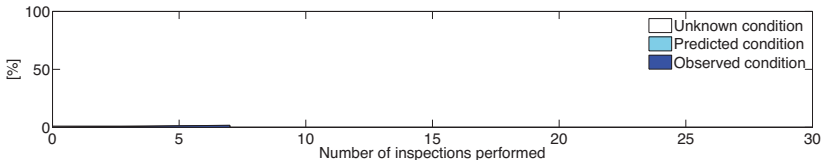


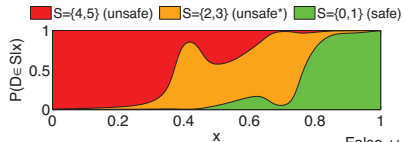
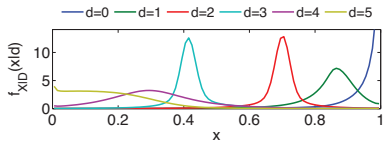
False +: 0% | False -: 0%



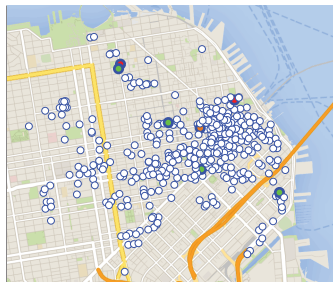
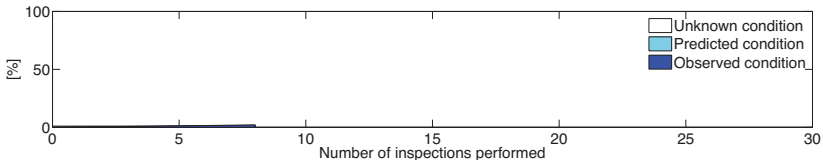


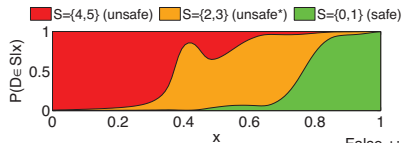
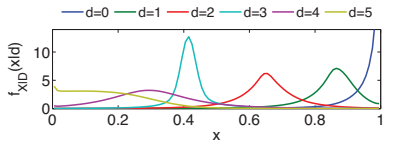
False +: 0% | False -: 0%



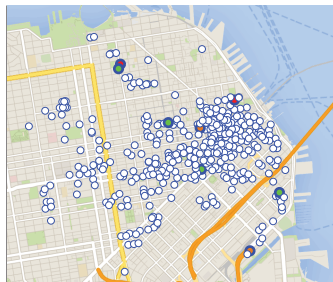
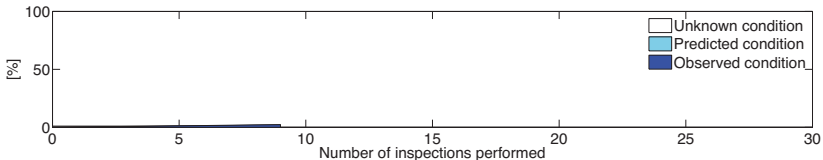


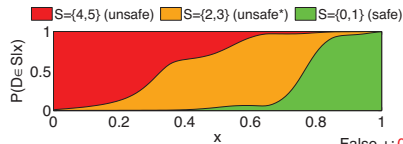
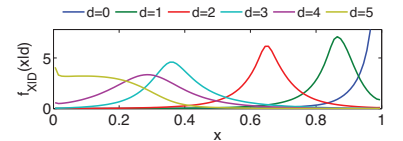
False +: 0% | False -: 0%



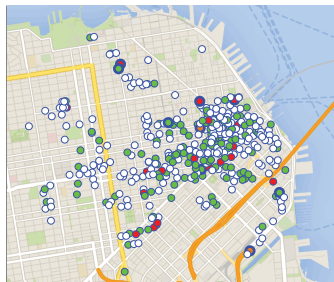
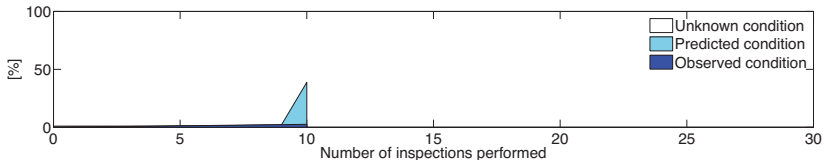


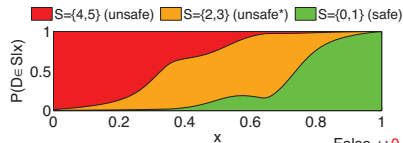
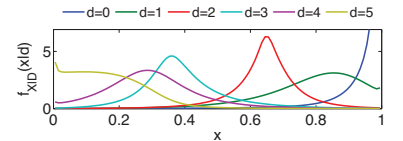
False +: 0% | False -: 0%



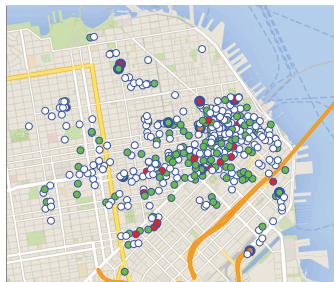
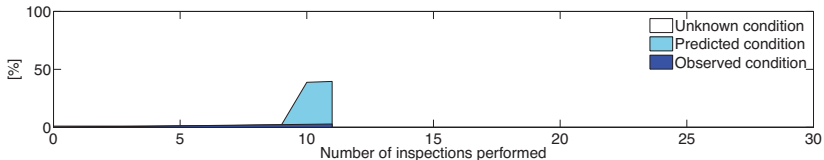


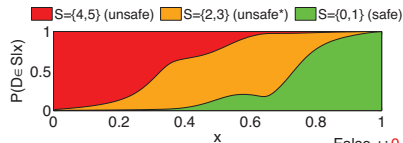
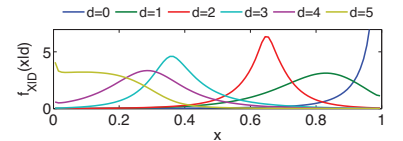
False +: 0% | False -: 3.9%



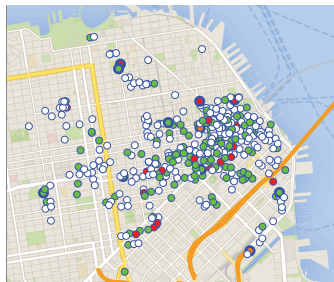
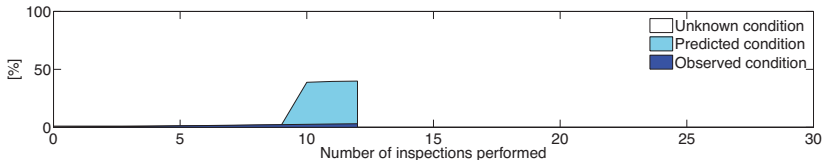


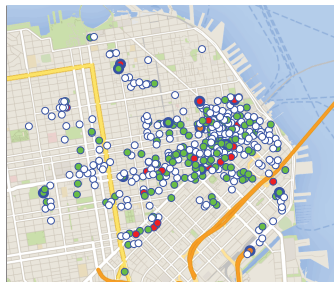
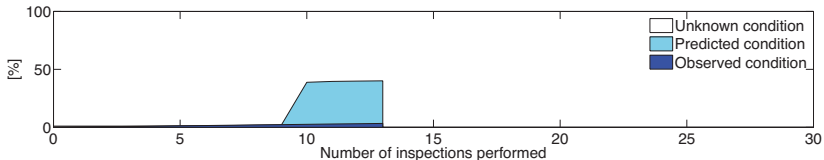
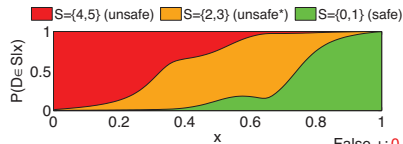
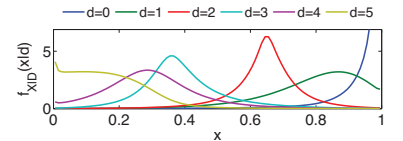
False +: 0.3% | False -: 3.9%

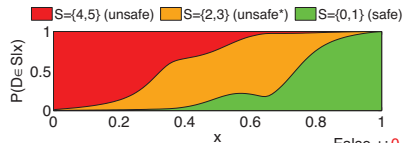
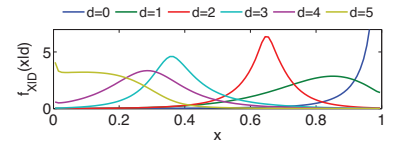




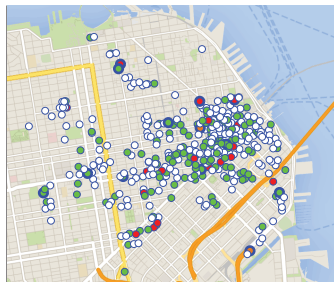
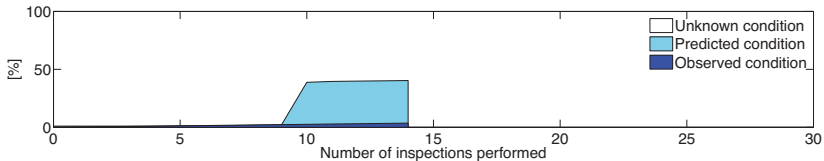
False +: 0.3% | False -: 3.9%

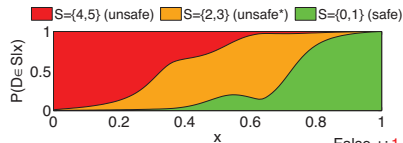
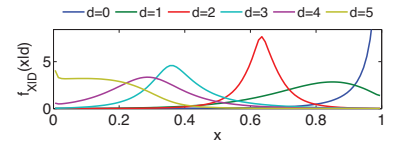




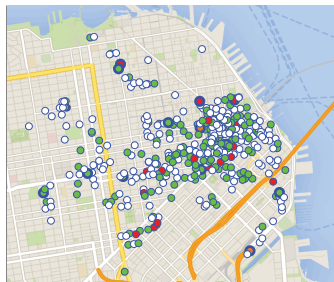
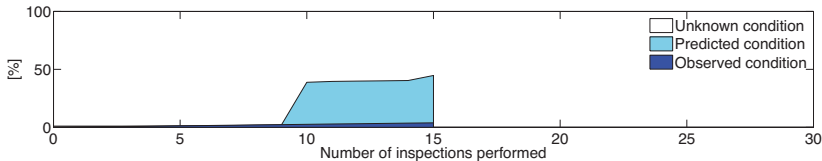


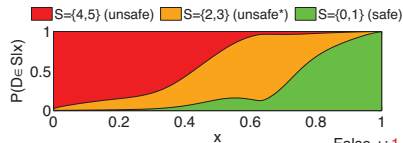
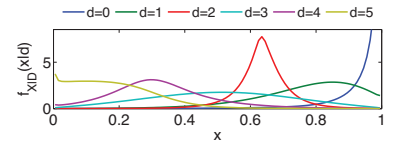
False +: 0.3% | False -: 3.9%



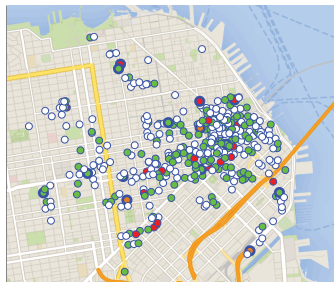
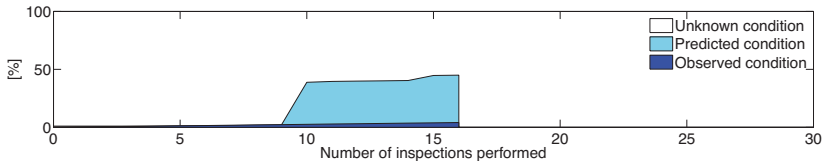


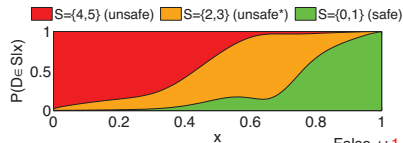
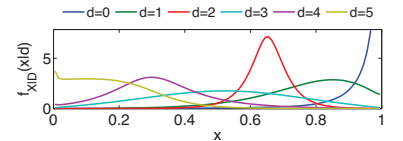
False +: 1.3% | False -: 3.9%



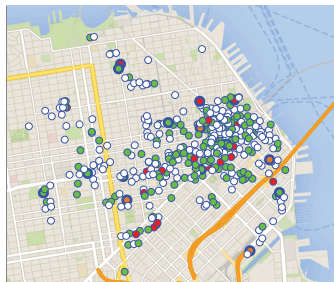
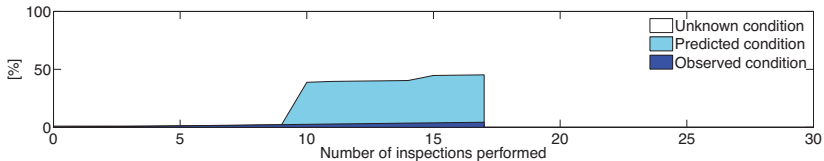


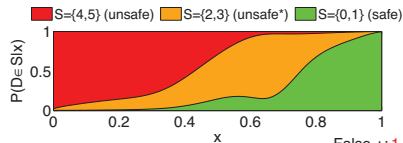
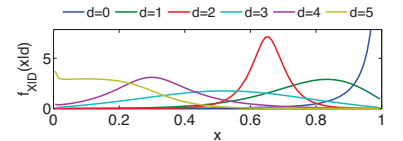
False +: 1.3% | False -: 3.9%



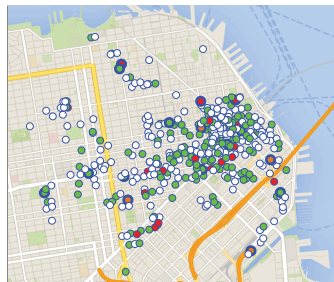
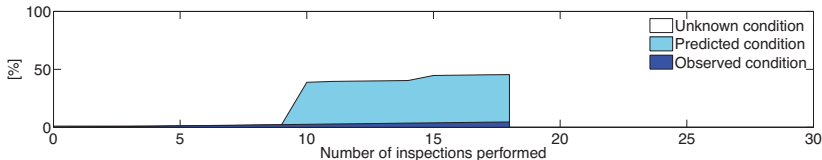


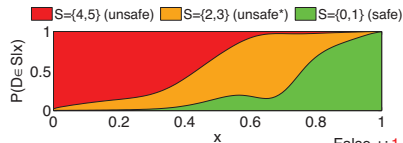
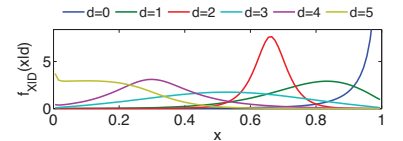
False +: 1.3% | False -: 3.9%



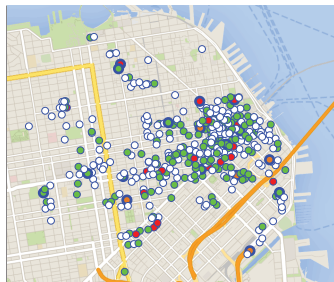
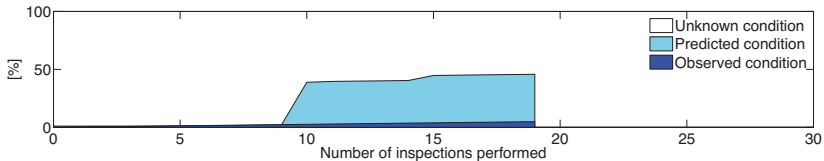


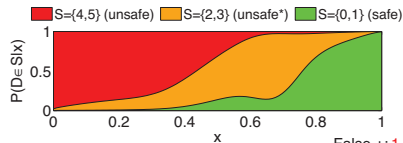
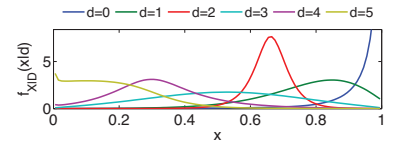
False +: 1.3% | False -: 3.9%



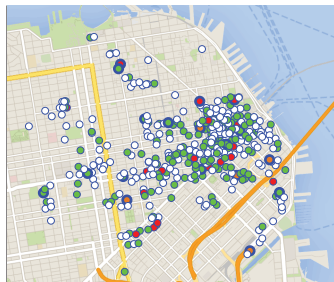
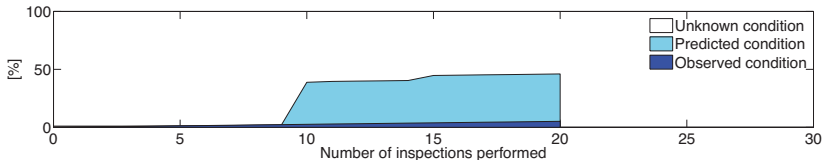


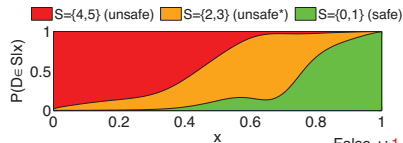
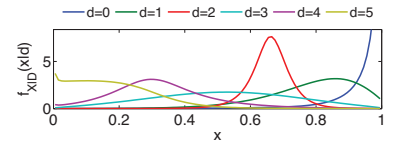
False +: 1.3% | False -: 3.9%



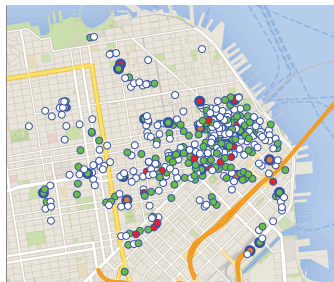
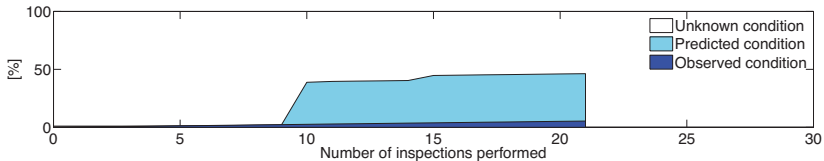


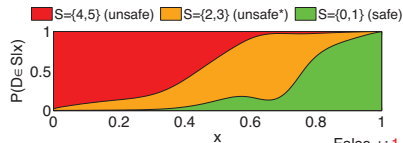
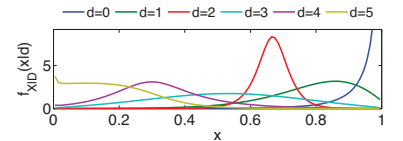
False +: 1.3% | False -: 3.9%



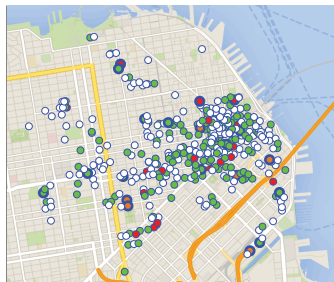
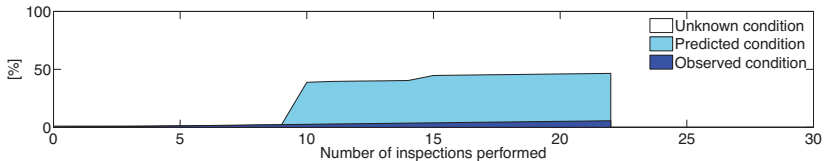


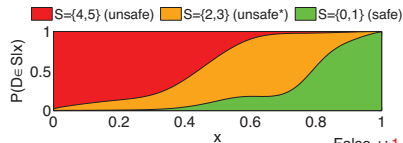
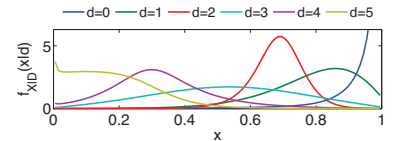
False +: 1.3% | False -: 3.9%



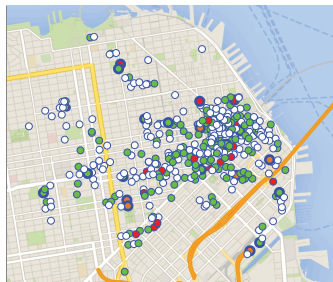
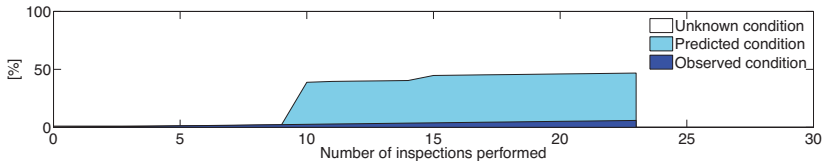


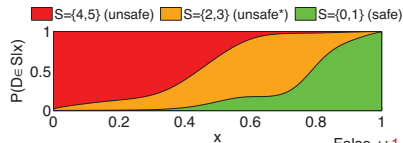
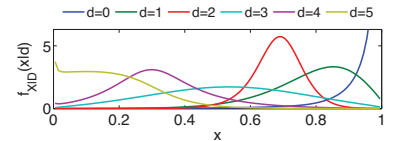
False +: 1.3% | False -: 3.9%



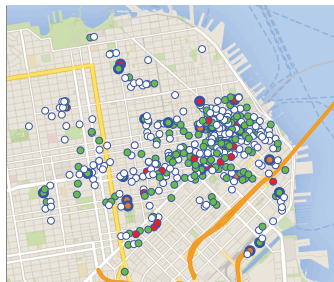
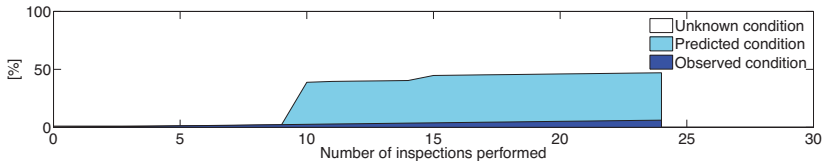


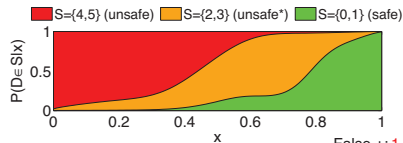
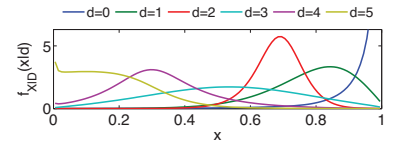
False +: 1.3% | False -: 3.9%



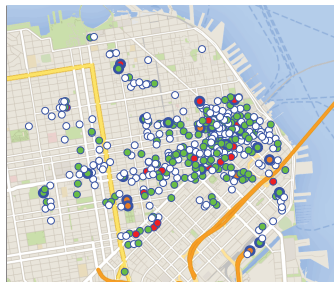
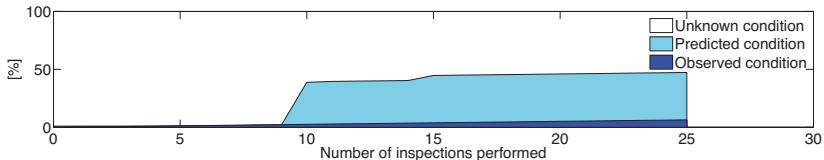


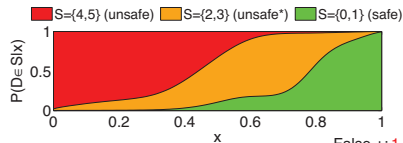
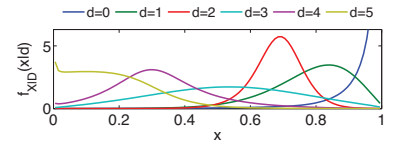
False +: 1.3% | False -: 3.9%



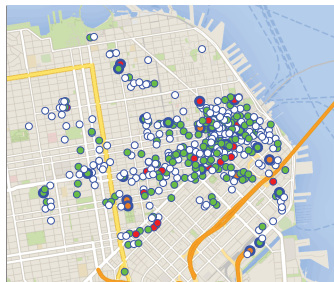
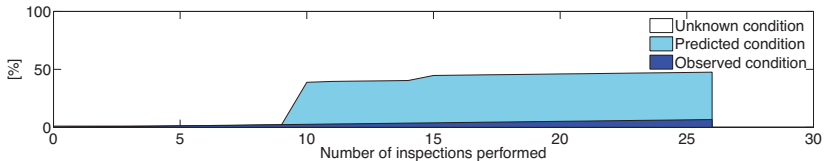


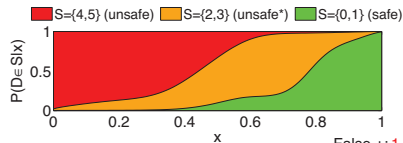
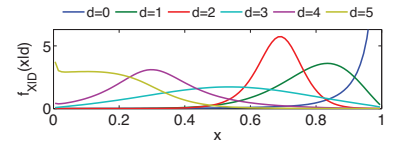
False +: 1.3% | False -: 3.9%



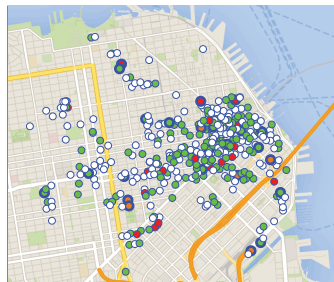
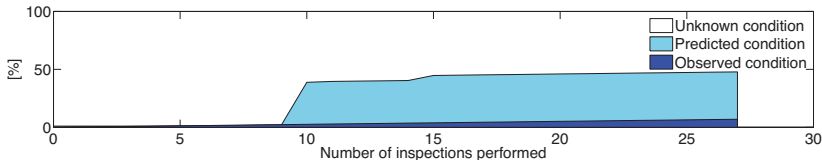


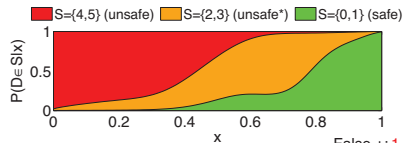
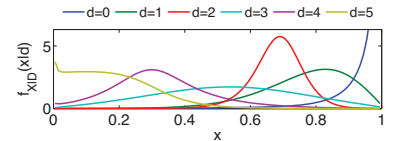
False +: 1.3% | False -: 3.9%



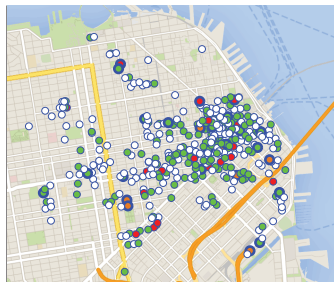
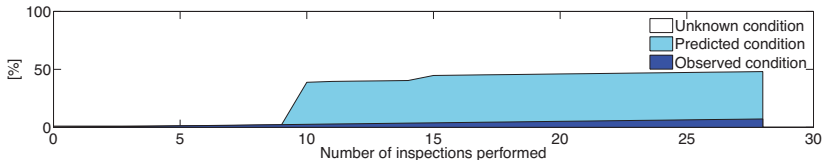


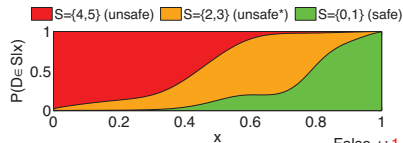
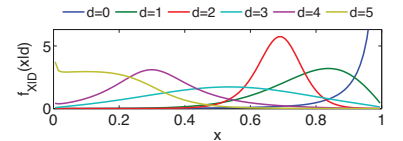
False +: 1.3% | False -: 3.9%



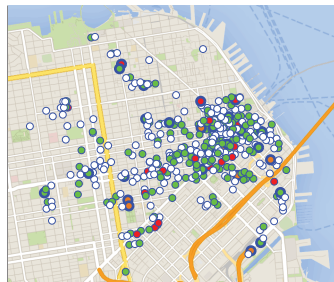
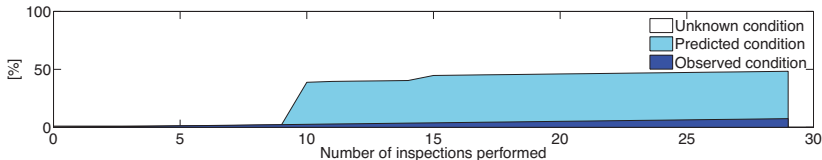


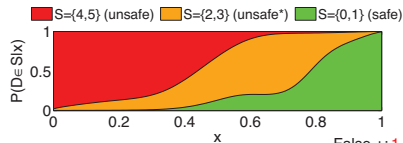
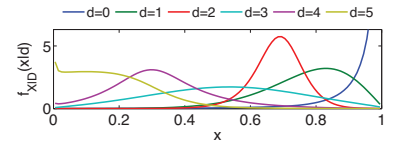
False +: 1.3% | False -: 3.9%



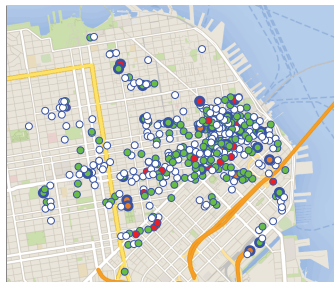
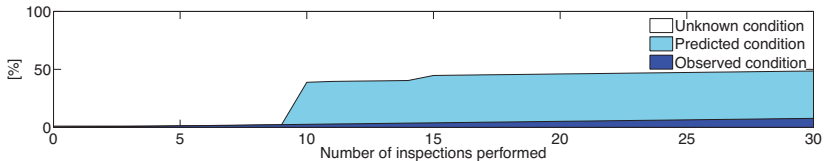


False +: 1.3% | False -: 3.9%

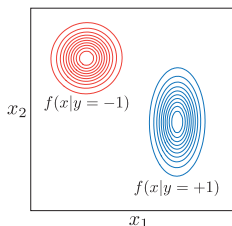




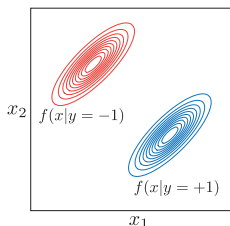
False +: 1.3% | False -: 3.9%



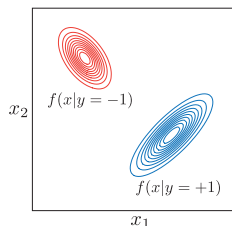
Classic generative methods



(a) Naive Bayes



(b) LDA



(c) QDA

Method	$f(\mathbf{x} y = j)$
Naive Bayes	$X_i \perp\!\!\!\perp X_k y, \forall i \neq k$ e.g. $\mathcal{B}(x_1; \alpha, \beta) \cdot \mathcal{N}(x_2, \mu, \sigma)$
LDA	$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{\mathbf{x},j}, \boldsymbol{\Sigma}_{\mathbf{x}})$
QDA	$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{\mathbf{x},j}, \boldsymbol{\Sigma}_{\mathbf{x},j})$

Summary - Generative classifier

MLE/MAP point estimation

$$\underbrace{p(y|x, \theta^*)}_{\text{posterior}} = \frac{\underbrace{f(x|y, \theta^*)}_{\text{likelihood}} \cdot \underbrace{p(y; \theta^*)}_{\text{prior}}}{\underbrace{f(x; \theta^*)}_{\text{norm. cte.}}}$$

Bayesian estimation

$$\underbrace{p(y|x, \mathcal{D})}_{\text{posterior predictive}} = \int \frac{\underbrace{f(x|y, \theta)}_{\text{likelihood}} \cdot \underbrace{p(y; \theta)}_{\text{prior}}}{\underbrace{f(x; \theta)}_{\text{norm. cte.}}} \cdot \underbrace{f(\theta|\mathcal{D})}_{\text{posterior}} d\theta$$

Summary - Strength & Limitations

$$\underbrace{p(y|x)}_{\text{posterior}} = \frac{\underbrace{f(x|y)}_{\text{likelihood}} \cdot \underbrace{p(y)}_{\text{prior}}}{\underbrace{f(x)}_{\text{norm. cte.}}}$$

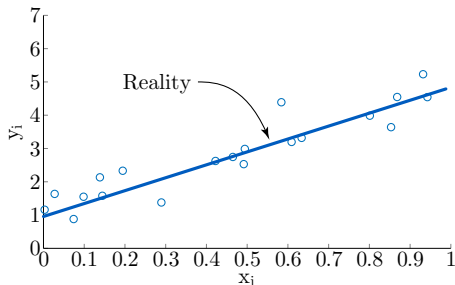
Strength

- ▶ Compatible with any type of observation (error type/observation type)
- ▶ Small dataset → Bayesian inference
- ▶ Can handle unbalanced observation classes

Limitations

- ▶ Complex to set-up tailored model structures
- ▶ Simple structures (e.g. Naive Bayes, LDA & QDA) are outperformed by other *discriminative* methods

Revision – Linear Regression [≡]

**Data**

$$\mathcal{D} = \{(x_i, y_i), \forall i = 1 : D\}$$

$$x_i \in \mathbb{R} : \begin{cases} \text{Covariate} \\ \text{attribute} \\ \text{regressor} \end{cases}$$

$$y_i \in \mathbb{R} : \text{Observation}$$

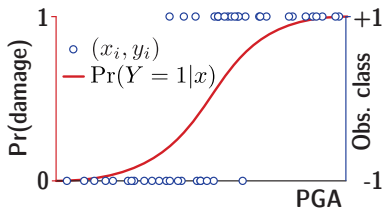
Model

$$g(x) \equiv \text{fct}(x)$$

Hypothesis: $g(x) \equiv \text{reality}$, i.e. no variability

$$y = g(x) + v, \quad v : V \sim \mathcal{N}(v; 0, \sigma_V^2)$$

Introduction to Logistic regression



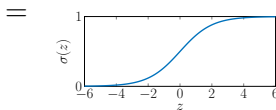
Data

$$\mathcal{D} = \{(x_i, y_i), \forall i = 1 : D\}$$

$$x_i \in \mathbb{R} : \begin{cases} \text{Covariate} \\ \text{attribute} \\ \text{regressor} \end{cases}$$

$$y_i \in \{-1, 1\} : \text{Observation}$$

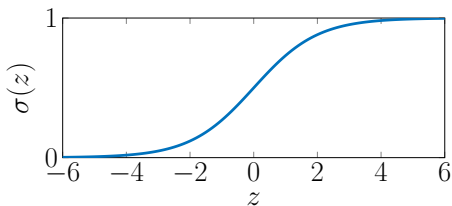
$$\begin{aligned} \Pr(Y = 1|x) &= \sigma(\overbrace{b_0 + b_1 x}^z) \\ &= \frac{1}{1 + \exp(-z)} \end{aligned}$$



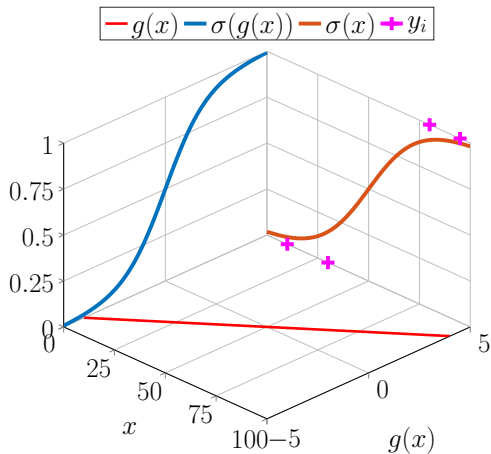
From linear to logistic regression

$$\underbrace{\mathbf{x} \in \mathbb{R}^X}_{\text{covariates}} \rightarrow \underbrace{g(\mathbf{x}) = \mathbf{X}\mathbf{b} \in \mathbb{R}}_{\text{hidden/latent variable}} \rightarrow \underbrace{\sigma(g(\mathbf{x})) \in (0, 1)}_{\text{Pr}(Y=y|\mathbf{x})}$$

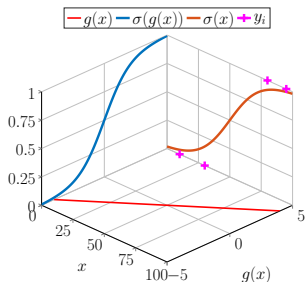
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$



Example - Linear regression & sigmoid []



Likelihood



Marginal likelihoods

For $x_i : y_i = +1 \rightarrow \mathbf{X}_{(+1)}$

$$\Pr(Y = +1|x_i) = \sigma(\mathbf{X}_{(+1)}\mathbf{b}) \equiv \sigma(x_i)$$

For $x_i : y_i = -1 \rightarrow \mathbf{X}_{(-1)}$

$$\Pr(Y = -1|x_i) = 1 - \sigma(\mathbf{X}_{(-1)}\mathbf{b})$$

Joint Log-likelihood

$$\ln p(\mathcal{D}|\mathbf{b}) = \sum \ln(\sigma(\mathbf{X}_{(+1)}\mathbf{b})) + \sum \ln(1 - \sigma(\mathbf{X}_{(-1)}\mathbf{b}))$$

Maximum Likelihood Estimation

$$\mathbf{b}^* = \arg \max_{\mathbf{b}} \ln p(\mathcal{D}|\mathbf{b})$$

Limitations

- ▶ ⚠ The performance depends on the capacity to **hand-pick the correct transformation functions**
(Difficult when $X > 1 \rightarrow$ **Cross-validation**)
- ▶ ⚠ No analytic formulation to estimate parameters
- ▶ ⚠ Only compatible with error-free direct observations
- ▶ ⚠ Not as powerful as modern methods

Why is it used then?

- ▶ ⚠ Simple
- ▶ ⚠ Interpretability of model parameters **b**
(Discrete choice models \rightarrow transportation & economics)

Section Outline

GPC

- 4.1 GPR v.s. GPC
- 4.2 Updating a GP using exact observations
- 4.3 GPC formulation

Gaussian Process Regression summary

Given a **system response** so that

$$\underbrace{g_i}_{\text{observation}} = \underbrace{g(\mathbf{x}_i)}_{\text{reality}} \in \mathbb{R}, \quad \underbrace{\mathbf{x}_i}_{\text{covariates}} = [x_1, x_2, \dots, x_X]^T \in \mathbb{R}^X$$

Data: $\mathcal{D} = \{(\mathbf{x}_i, g_i), \forall i = 1 : D\}$

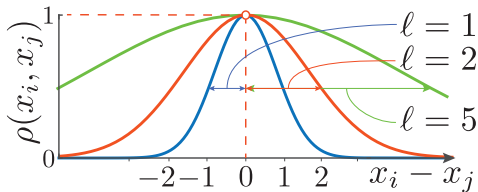
$$\mathbf{g} = [g_1, g_2, \dots, g_D]^T_{1 \times D} \in \mathbb{R}^D, \quad \mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_X]_{D \times X}$$

Gaussian process: $g(\mathbf{x}) : \mathbf{G} \sim \mathcal{N}(g(\mathbf{x}); \mathbf{m}_G, \Sigma_G)$

Set of discrete Gaussian random variables for which the pairwise correlation between G_i and G_j is a function of the distance between attributes

$$[\Sigma_G]_{ij} = \rho(x_i, x_j) \sigma_{G_i} \sigma_{G_j}, \quad \rho(x_i, x_j) = \text{fct}(|x_i - x_j|)$$

Correlation function



Updating a GP using exact observations

Given $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{g}_i), i = 1, \dots, D\}$ a set of D observations and \mathbf{x}_* a set of \mathbf{X}_* covariates for which we want to predict

$$f(\mathbf{g}_* | \mathbf{x}_*, \mathcal{D})$$

Reminder: Gaussian conditionals are also gaussian

$$\underbrace{\left\{ \begin{array}{c} \mathbf{G} \\ \mathbf{G}_* \end{array} \right\}, \mathbf{m} = \left\{ \frac{\mathbf{m}_G}{\mathbf{m}_{G_*}} \right\}, \boldsymbol{\Sigma} = \left[\begin{array}{c|c} \boldsymbol{\Sigma}_G & \boldsymbol{\Sigma}_{G_*} \\ \hline \boldsymbol{\Sigma}_{G_*}^\top & \boldsymbol{\Sigma}_* \end{array} \right]}_{\text{Prior knowledge}}$$

Prior knowledge

$$\begin{aligned} [\boldsymbol{\Sigma}_G]_{ij} &= \rho(x_i, x_j) \sigma_G^2 \\ [\boldsymbol{\Sigma}_{G_*}]_{ij} &= \rho(x_i, x_{*j}) \sigma_G^2 \\ [\boldsymbol{\Sigma}_*]_{ij} &= \rho(x_{*i}, x_{*j}) \sigma_G^2 \end{aligned}$$

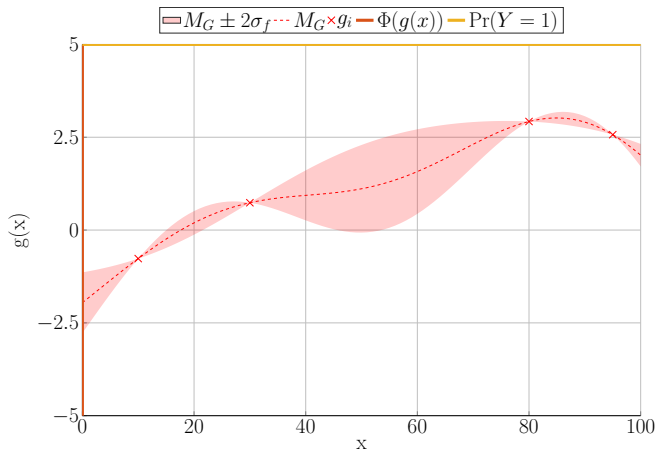
$$f(\mathbf{g}_* | \mathbf{x}_*, \mathcal{D}) = \mathcal{N}(\mathbf{g}_*; \mathbf{m}_{*|\mathcal{D}}, \boldsymbol{\Sigma}_{*|\mathcal{D}})$$

$$\mathbf{m}_{*|\mathcal{D}} = \mathbf{m}_{G_*} + \boldsymbol{\Sigma}_{G_*}^\top \boldsymbol{\Sigma}_G^{-1} (\mathbf{g} - \mathbf{m}_G)$$

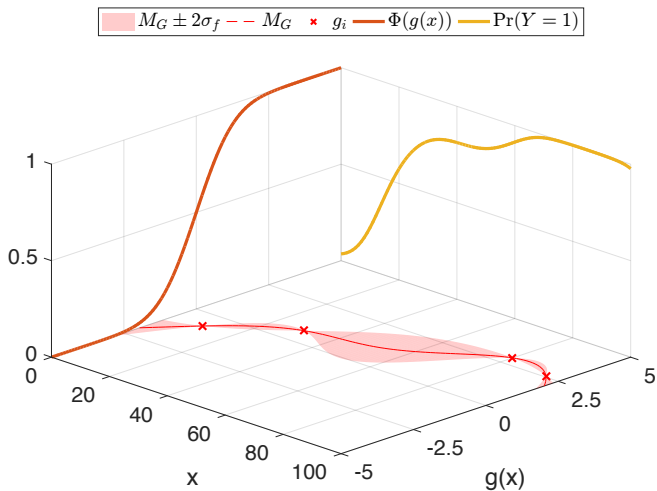
$$\boldsymbol{\Sigma}_{*|\mathcal{D}} = \boldsymbol{\Sigma}_* - \boldsymbol{\Sigma}_{G_*}^\top \boldsymbol{\Sigma}_G^{-1} \boldsymbol{\Sigma}_{G_*}$$

Posterior knowledge

Example - GPR v.s. GPC [CIV_ML]

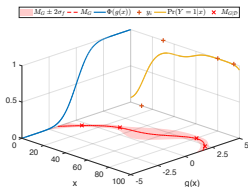


Example - GPR v.s. GPC [CIV_ML]



[CIV_ML/GPC_exemple_1.m]

Example - GPR v.s. GPC [CIV_ML]



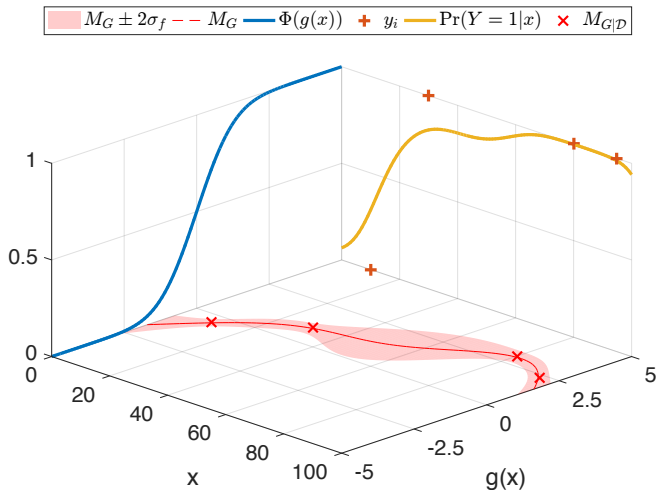
$\Phi(\cdot) \triangleq$ Normal CDF

$$\begin{aligned} \Pr(Y = 1 | x_*, \mathcal{D}) &= \int \Phi(g_*) f(g_* | x_*, \mathcal{D}) dg_* \\ &= \Phi\left(\frac{m_{G_*|g}}{\sqrt{1 + \text{var}[G_*|g]}}\right) \end{aligned}$$

**In a classification setup, g is not observed
Only $y_i \in \{-1, 1\}$ is observed.**

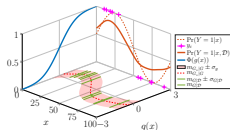
For each $y_i \in \{-1, 1\}$ we need to infer g_i

Example - GPR v.s. GPC [CIV_ML/GPC_exemple_2.m]



Gaussian Process Classification Formulation

$$\begin{aligned} \Pr(Y = 1|x_*, \mathcal{D}) &= \int \Phi(g_*) \cdot f(g_*|\mathbf{g}, x_*, \mathcal{D}) dg_* \\ &= \Phi\left(\frac{\mathbb{E}[G_*|x_*, \mathbf{g}, \mathcal{D}]}{\sqrt{1+\text{var}[G_*|x_*, \mathbf{g}, \mathcal{D}]}}\right) \end{aligned}$$



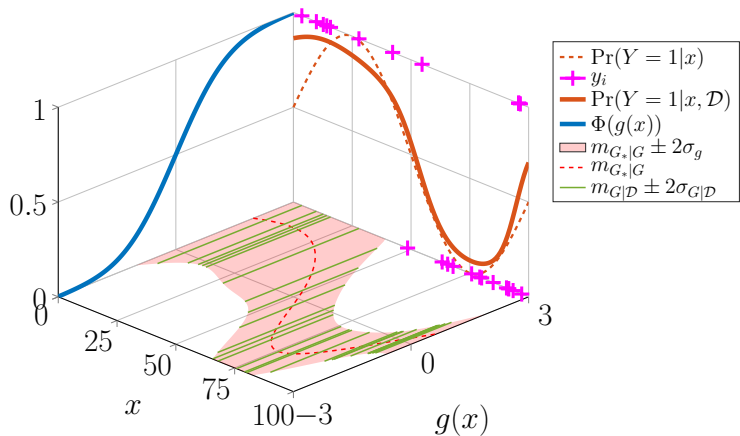
$$f(\mathbf{g}_*|\mathbf{g}) \equiv \mathcal{N}(\mathbf{g}_*; \boldsymbol{\mu}_{*|\mathbf{g}, \mathcal{D}}, \boldsymbol{\Sigma}_{*|\mathbf{g}, \mathcal{D}}) \quad (\text{GPR})$$

$$\boldsymbol{\mu}_{*|\mathbf{g}, \mathcal{D}} = \boldsymbol{\Sigma}_{\mathbf{G}_*}^T \boldsymbol{\Sigma}_{\mathbf{G}}^{-1} (\boldsymbol{\mu}_{\mathbf{G}|\mathcal{D}} - \boldsymbol{\mu}_{\mathbf{G}} - \mathbf{0}) \quad (\text{conditional mean})$$

$$\boldsymbol{\Sigma}_{*|\mathbf{g}, \mathcal{D}} = \boldsymbol{\Sigma}_* - \boldsymbol{\Sigma}_{\mathbf{G}_*}^T \boldsymbol{\Sigma}_{\mathbf{G}|\mathcal{D}}^{-1} \boldsymbol{\Sigma}_{\mathbf{G}_*} \quad (\text{conditional COV})$$

$$\begin{aligned} f(\mathbf{g}|\mathcal{D}) &= \frac{p(\mathcal{D}_y|\mathbf{g}) \cdot f(\mathbf{g}|\mathcal{D}_x)}{p(\mathcal{D}_y|\mathcal{D}_x)} && (\text{posterior}) \\ &\propto p(\mathcal{D}_y|\mathbf{g}) \cdot f(\mathbf{g}|\mathcal{D}_x) && (\text{hidden covariates}) \\ &\approx \mathcal{N}(\mathbf{g}; \boldsymbol{\mu}_{\mathbf{G}|\mathcal{D}}, \boldsymbol{\Sigma}_{\mathbf{G}|\mathcal{D}}) && (\text{Laplace approx.}) \end{aligned}$$

Example - GPC / GPML



Section Outline

Neural Networks

- 5.1 Structure - From Logistic regression to NN
 - 5.2 Structure - Bivariate classification
 - 5.3 Structure - Multivariate classification
 - 5.4 Deep Learning – CNN
 - 5.5 Examples
-

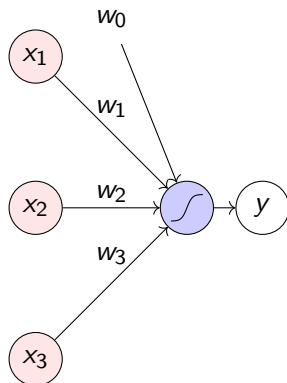
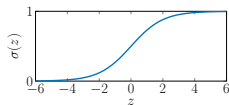
Introduction to the neural network structure

Data $\mathcal{D} = \{(x_i, y_i), \forall i = 1 : D\}$

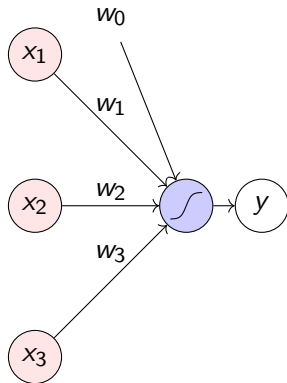
$x_i \in \mathbb{R} :$ $\begin{cases} \text{Covariate} \\ \text{attribute} \\ \text{regressor} \end{cases}$

$y_i \in \{-1, 1\} :$ Observation

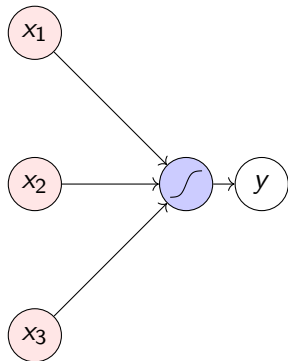
$$\begin{aligned} Pr(Y = 1|\mathbf{x}) &= \sigma(\overbrace{w_0 + w_1x_1 + w_2x_2 + w_3x_3}^z) \\ &= \frac{1}{1 + \exp(-z)} \\ &= \end{aligned}$$



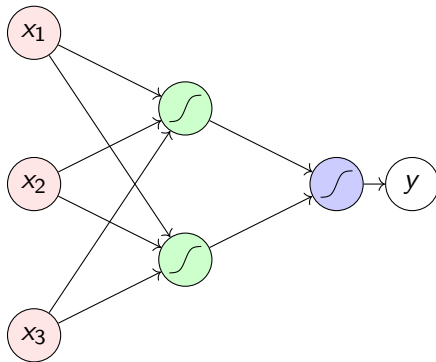
Neural network structure for binary classification



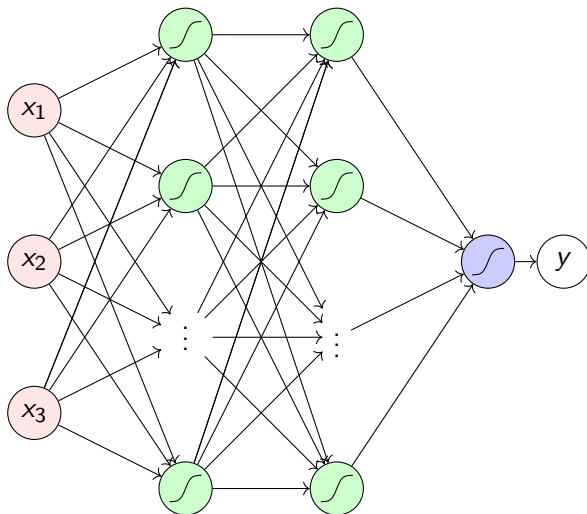
Neural network structure for binary classification



Neural network structure for binary classification



Neural network structure for binary classification



NN – Multivariate regression

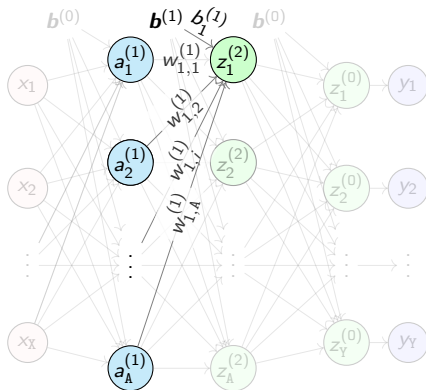
\mathbf{x} : Input layer / covariates

$\mathbf{a}^{(l)}$: Activation units
 $a_i^{(j)} = \sigma(z_i^{(j)})$

$\mathbf{z}^{(l)}$: Hidden units
 $z_i^{(2)} = \mathbf{w}^T \mathbf{a}^{(1)} + b_i$

$\mathbf{z}^{(0)}$: Output layer

\mathbf{y} : Observations



NN – Multivariate classification

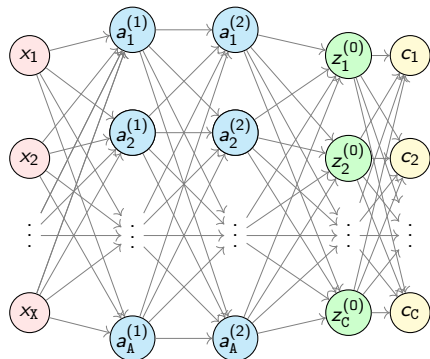
\mathbf{x} : Input layer / covariates

$\mathbf{a}^{(l)}$: Activation units

$\mathbf{z}^{(l)}$: Hidden units

$\mathbf{z}^{(0)}$: Output layer, $\mathbf{z}^{(0)} \in \mathbb{R}^C$

\mathbf{c} : Classes, $\mathbf{c} \in \{0, 1\}^C$

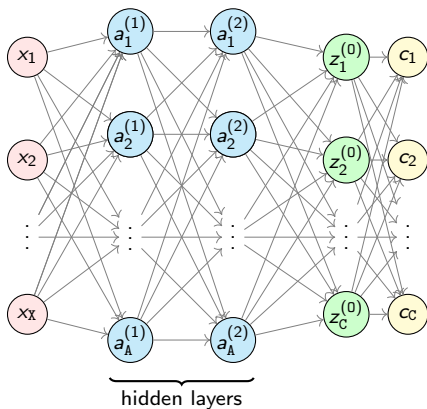


$$\Pr(c_i) = \text{softmax}(\mathbf{z}^{(0)}, i) = \frac{\exp(z_i^{(0)})}{\sum_{j=1}^C \exp(z_j^{(0)})}$$

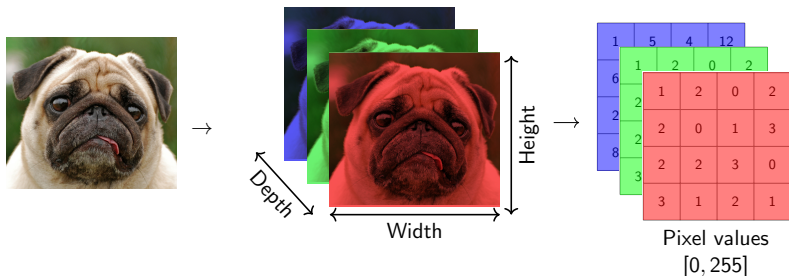
Deep Learning

Deep learning refers to a neural network which has a deep structure

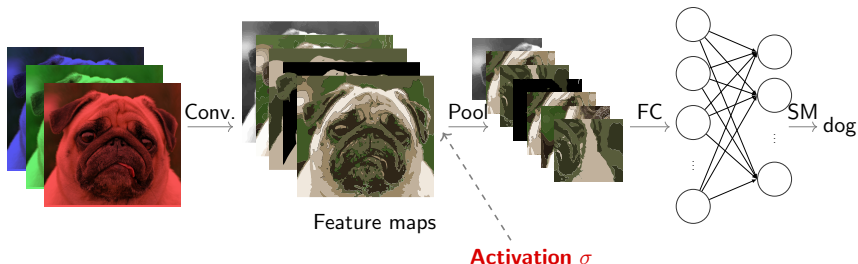
i.e., a large number of **hidden layers**



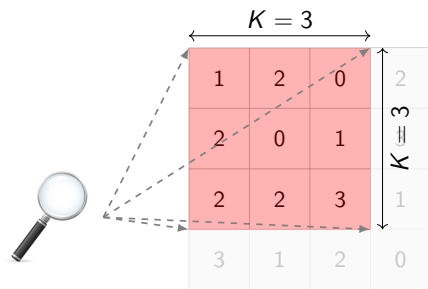
Convolutional Neural Networks (CNNs)



Convolutional Neural Networks (CNNs)



Kernel Size (K) & Stride (S)



Zero-padding (P)

1	2	0	2		
2	0	1	3		
2	2	3	1		
3	1	2	0		

$\xrightarrow{P=1}$

0	0	0	0	0	0
0	1	2	0	2	0
0	2	0	1	3	0
0	2	2	3	1	0
0	3	1	2	0	0
0	0	0	0	0	0

$$W^0 = \frac{W^I - K + 2P}{S} + 1$$

$$H^0 = \frac{H^I - K + 2P}{S} + 1$$

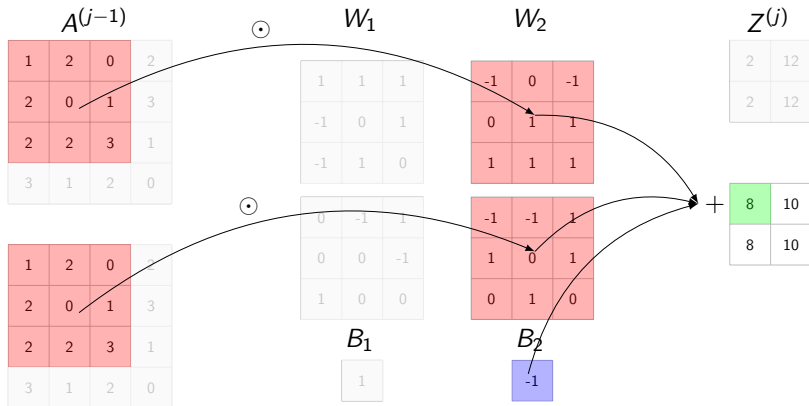
Given $K = 3$, $S = 1$, $P = 1$

$$W^0 = \frac{4 - 3 + 2 \times 1}{1} + 1$$

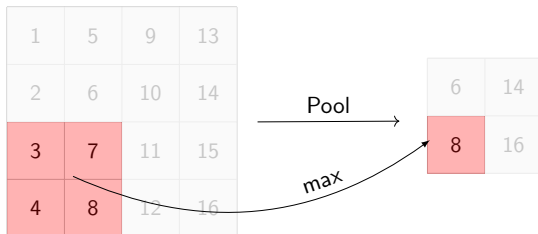
$$= 4$$

$$H^0 = 4$$

Convolutional Layers – $\{K = 3, S = 1, P = 0\}$



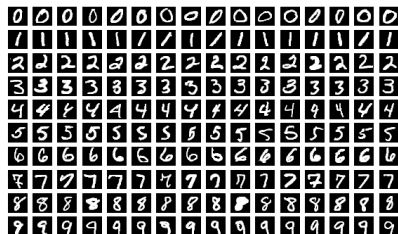
Pooling Layers – $\{K = 2, S = 2\}$



Max Pooling

MNIST

train + test 784
 60 000 + 10 000 28 × 28 images



Network's Architecture

(Wan et al., 2013)

Layer	$D \times W \times H$	$K \times K$	P	S	σ
Input	$1 \times 28 \times 28$	-	-	-	-
Conv.	$32 \times 27 \times 27$	4×4	1	1	ReLU
Pooling	$32 \times 13 \times 13$	3×3	0	2	-
Conv.	$64 \times 9 \times 9$	5×5	0	1	ReLU
Pooling	$64 \times 4 \times 4$	3×3	0	2	-
FC	$150 \times 1 \times 1$	-	-	-	ReLU
Output	$10 \times 1 \times 1$	-	-	-	-

	Error Rate [%] Hyperparameters	
	$E = e$	$e \quad B$
TAGI	0.56	13 1
Backprop	0.67	1000 128

[Nguyen and Goulet (2020)]

⌘ Polytechnique Montréal

CIFAR10



Network's Architecture

(Wan et al., 2013)

Layer	$D \times W \times H$	$K \times K$	P	S	σ
Input	$3 \times 32 \times 32$	-	-	-	-
Conv.	$32 \times 32 \times 32$	5×5	2	1	ReLU
Pooling	$32 \times 16 \times 16$	3×3	1	2	-
Conv.	$32 \times 16 \times 16$	5×5	2	1	ReLU
Pooling	$32 \times 8 \times 8$	3×3	1	2	-
Conv.	$64 \times 8 \times 8$	5×5	2	1	ReLU
Pooling	$64 \times 4 \times 4$	3×3	1	2	-
FC	$64 \times 1 \times 1$	-	-	-	ReLU
Output	$10 \times 1 \times 1$	-	-	-	-

	Error Rate [%] Hyperparameters	
	$E = e$	e B
TAGI	23.6	50 16
Backprop	23.5	150 128

[pytorch.com]

Summary

Generative Classifier:

- ▶ Offers the most flexibility (e.g. indirect observations)
- ▶ Best suited for small datasets

Gaussian Process Classification:

- ▶ Quick & easy
- ▶ Allows interpolating and extrapolating (uncertainty estimates)
- ▶ Best suited for medium size datasets

Neural networks / Deep Learning:

- ▶ State-of-the-art method
- ▶ Best suited for large datasets

Temporary page!

\LaTeX was unable to guess the total number of pages correctly. As there was some unprocessed data that should have been added to the final page this extra page has been added to receive it. If you rerun the document (without altering it) this surplus page will go away, because \LaTeX now knows how many pages to expect for this document.