

MITx: Statistics, Computation & Applications

Genomics and High-Dimensional Data Module

Lecture 2: Classification with High-Dimensional Data

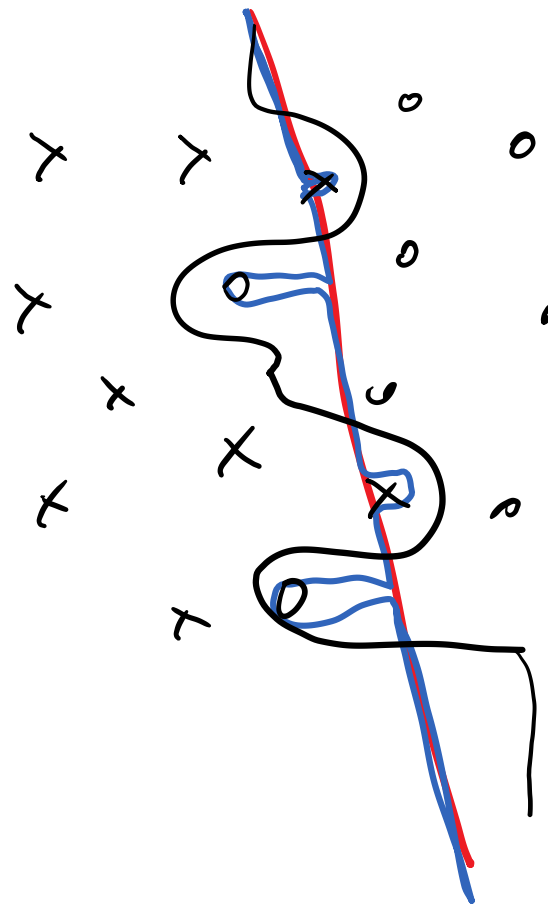
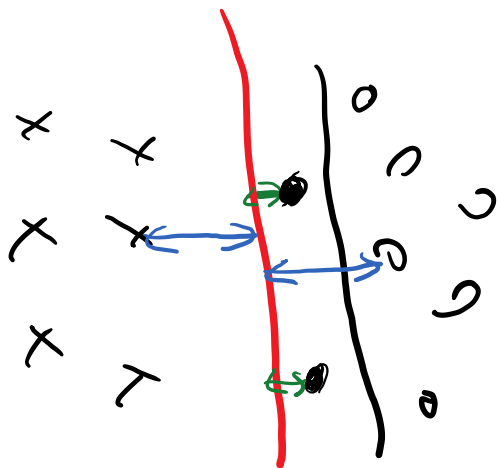
$x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^p$
 $c^{(1)}, \dots, c^{(n)} \in \mathcal{C}$

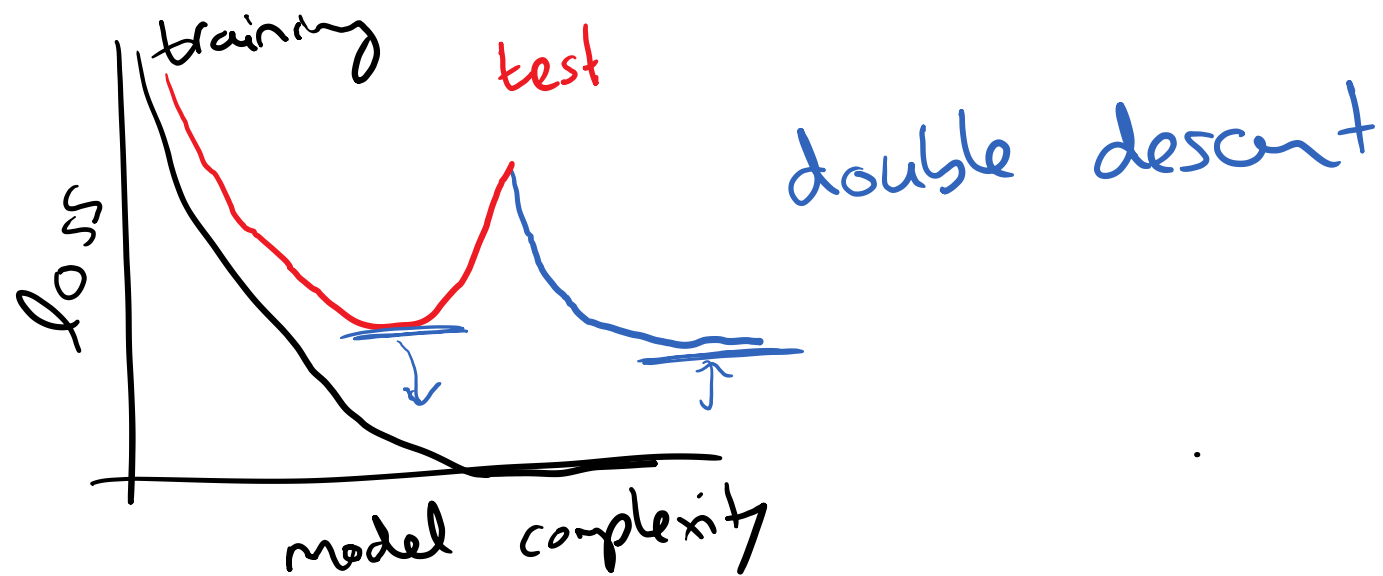
$f: \mathbb{R}^p \rightarrow \mathcal{C}$
 class labels

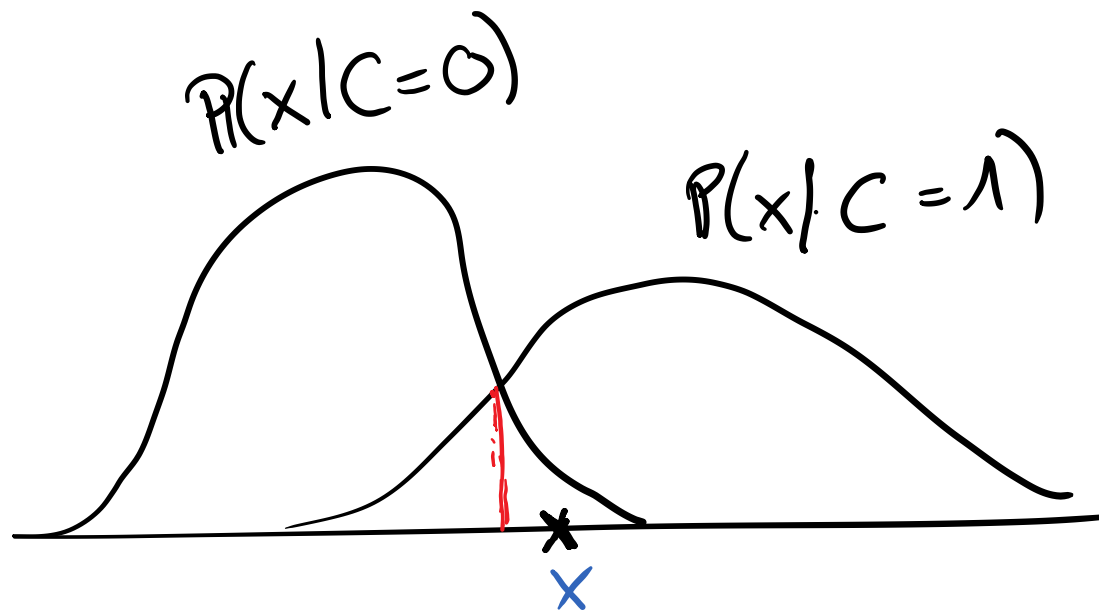
Confusion matrix
 estimated class

true class	estimated class			
	1	2	3	4
1	m			
2		n		
3			k	
4				m

SVM







$$P(C=0) \gg P(C=1)?$$

Best guess $x \in C=1$

$$\left. \begin{array}{l} \text{Estimate } P(x|C) \sim \mathcal{N}(\mu_c, \Sigma_c) \\ P(C) \end{array} \right\} \text{QDA}$$

$$\begin{array}{ll}
 P(C=0|x) & P(C=1|x) \\
 \overset{\cap}{[0,1]} \hookrightarrow \beta_0 + \beta^T x & \in (-\infty, \infty)
 \end{array}$$

$$\log \frac{P(C=0|x)}{P(C=1|x)} = \beta_0 + \beta^T x + \lambda \|\beta\|_1 \implies \beta \text{ sparse}$$

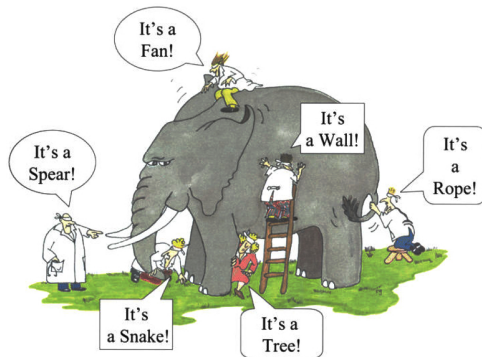
MITx: Statistics, Computation & Applications

Genomics and High-Dimensional Data Module

Lecture 2: Classification with Hig-Dimensional Data

Classification with high-dimensional data

- Linear / Quadratic discriminant analysis
- Logistic regression
- Support Vector Machines

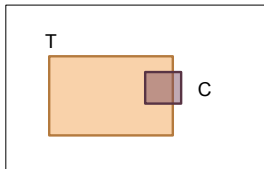


Refresher: Bayes rule

T: Med. Test positive

C: Patient has cancer

Sample space



(Marginal) Probability:
 $P(T)$, $P(C)$

New sample space:
People with cancer

$P(T|C)$
large



Conditional Probability:
 $P(T|C)$, $P(C|T)$

New sample space:
People with pos. test



$P(C|T)$
small

Bayes Theorem:

$$\text{posterior} \rightarrow P(C|T) = \frac{P(T|C)P(C)}{P(T)} \leftarrow \text{prior}$$

Class conditional probability

Using Bayes rule for classification

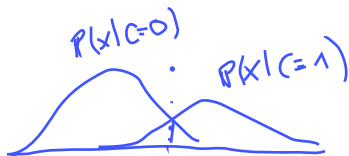
$$P(C|X) = \frac{P(C)P(X|C)}{P(X)} \sim P(C)P(X|C)$$

Find some estimate

Prior / prevalence:
Fraction of samples
in that class

Assume:
 $X|C \sim N(\mu_c, \Sigma_c)$

- Choose class $c \in \{1, \dots, K\}$ such that $P(C = c | X)$ is maximal



Using Bayes rule for classification

$$P(C|X) = \frac{P(C)P(X|C)}{P(X)} \sim P(C)P(X|C)$$

Find some estimate

Prior / prevalence:
Fraction of samples
in that class

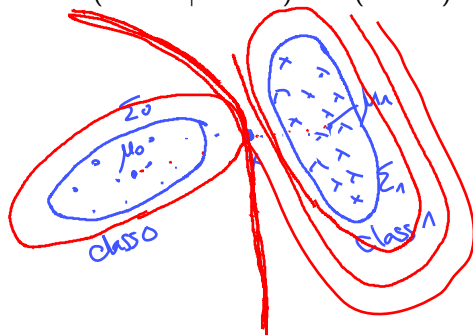
Assume:
 $X|C \sim \underline{N(\mu_c, \Sigma_c)}$

- Choose class $c \in \{1, \dots, K\}$ such that $P(C = c | X)$ is maximal
- Special case: 2 classes 0/1
 - choose $c = 1$ if $P(C = 1 | X) > 0.5$
 - equivalently, choose $c = 1$ if posterior odds $P(C = 1 | X) / P(C = 0 | X) > 1$
- We need to estimate $P(C = c)$ and $P(X | C)$

estimate:
 μ_c, Σ_c MLE

Quadratic discriminant analysis

- Assume $X \mid C = c \sim \mathcal{N}(\mu_c, \Sigma_c)$
- Estimate $P(C = c)$, μ_c , and Σ_c for each c (How?) MLE
- Choose class c such that $P(C = c \mid X = x) \propto P(C = c)P(X = x \mid C = c)$ is maximal



\Rightarrow classify x as class 0.

Quadratic discriminant analysis

- Assume $X \mid C = c \sim \mathcal{N}(\mu_c, \Sigma_c)$
- Estimate $P(C = c)$, μ_c , and Σ_c for each c (How?)
- Choose class c such that $P(C = c \mid X = x) \propto P(C = c)P(X = x \mid C = c)$ is maximal
- Use the fact that maximizing $P(C = c \mid X = x)$ is equivalent to maximizing $\log(P(C = c \mid X = x))$
- Do the math:
$$\log(P(C = c \mid X = x)) \propto \log(P(C = c)) - \frac{1}{2} \log \det \Sigma_c - \frac{1}{2} (x - \mu_c)^T \Sigma_c^{-1} (x - \mu_c)$$

$$f_{\mu_c, \Sigma_c} \propto \frac{1}{\det(\Sigma_c)^{1/2}} \exp\left(-\frac{1}{2} (x - \mu_c)^T \Sigma_c^{-1} (x - \mu_c)\right)$$

Quadratic discriminant analysis

- Assume $X \mid C = c \sim \mathcal{N}(\mu_c, \Sigma_c)$
- Estimate $P(C = c)$, μ_c , and Σ_c for each c (How?)
- Choose class c such that $P(C = c \mid X = x) \propto P(C = c)P(X = x \mid C = c)$ is maximal
- Use the fact that maximizing $P(C = c \mid X = x)$ is equivalent to maximizing $\log(P(C = c \mid X = x))$
- Do the math:
$$\log(P(C = c \mid X = x)) \propto \log(P(C = c)) - \frac{1}{2} \log \det \Sigma_c - \frac{1}{2} (\underline{x} - \underline{\mu}_c)^T \underline{\Sigma}_c^{-1} (\underline{x} - \underline{\mu}_c)$$
- Decision boundaries are quadratic $P(C=0 \mid X=x) = P(C=1 \mid X=x)$

Linear discriminant analysis

- Assume same covariance matrix Σ in all classes, i.e.
 $X \mid C = c \sim \mathcal{N}(\mu_c, \underline{\Sigma})$
- Estimate $P(C = c)$, μ_c , and Σ for each c
- Choose class c such that $\log(P(C = c \mid X = x))$ is maximal
- Do the math:

$$\log(P(C = c \mid X = x)) \propto \log(P(C = c)) - \frac{1}{2}\mu_c^T \Sigma^{-1} \mu_c + x^T \Sigma^{-1} \mu_c$$

- Decision boundaries are **linear**



Linear discriminant analysis

- Assume same covariance matrix Σ in all classes, i.e.
 $X \mid C = c \sim \mathcal{N}(\mu_c, \Sigma)$
- Estimate $P(C = c)$, μ_c , and Σ for each c
- Choose class c such that $\log(P(C = c \mid X = x))$ is maximal
- Do the math:

$$\log(P(C = c \mid X = x)) \propto \log(P(C = c)) - \frac{1}{2}\mu_c^T \Sigma^{-1} \mu_c + x^T \Sigma^{-1} \mu_c$$

- Decision boundaries are **linear**

What happens if we assume $\Sigma = \sigma^2 I$?



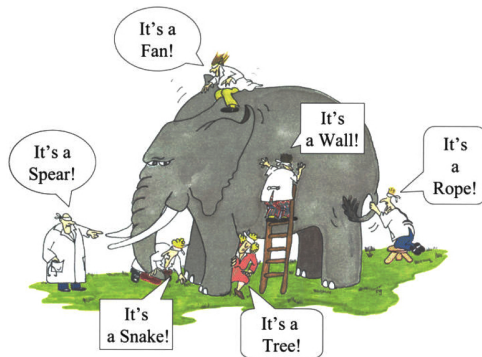
$$X^T X = \Sigma = U \Lambda U^T$$

$$\Lambda^{-1/2} U^T (X^T X) U \Lambda^{-1/2} = \Lambda^{-1/2} U^T (U \Lambda U^T) U \Lambda^{-1/2} = \text{Id}$$

$$X \rightarrow X U \Lambda^{-1/2}$$

Classification with high-dimensional data

- Linear / Quadratic discriminant analysis
- Logistic regression
- Support Vector Machines



Reduced-rank LDA (a.k.a. Fisher's LDA)

How can we use LDA to find informative low-dimensional projection of the data?

Reduced-rank LDA (a.k.a. Fisher's LDA)

How can we use LDA to find informative low-dimensional projection of the data?

- idea: $\mu_1, \dots, \mu_K \in \mathbb{R}^p$ lie in a linear subspace of $\dim K - 1$ (usually $p \gg k$)
- if $K = 3$, then data can be projected into 2d
- if $K > 3$, combine LDA with PCA, i.e. perform PCA on class means
 - 1. LD is 1. PC of class means, 2. LD is 2. PC of class means, etc.

Reduced-rank LDA (a.k.a. Fisher's LDA)

How can we use LDA to find informative low-dimensional projection of the data?

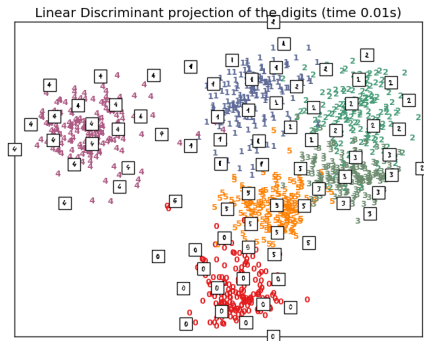
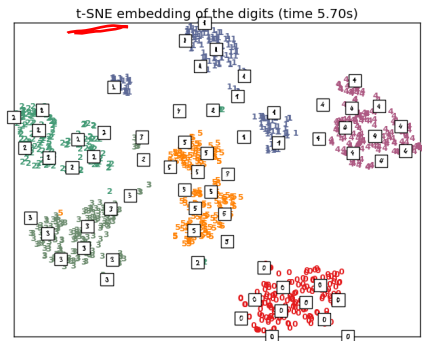
- idea: $\mu_1, \dots, \mu_K \in \mathbb{R}^p$ lie in a linear subspace of $\dim K - 1$ (usually $p \gg k$)
- if $K = 3$, then data can be projected into 2d
- if $K > 3$, combine LDA with PCA, i.e. perform PCA on class means
 - 1. LD is 1. PC of class means, 2. LD is 2. PC of class means, etc.

What is maximum number of LDs?

$$\min(p, \underline{k-1})$$

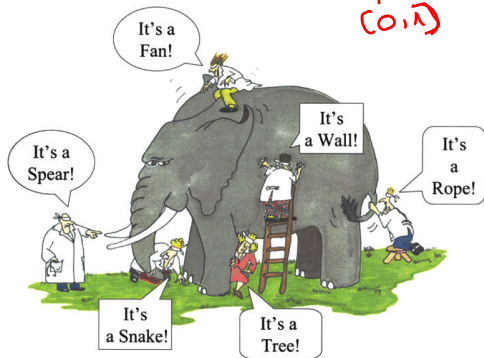
Example: Digit recognition

LDA can be used to find informative low-dimensional projection of the data by performing PCA on class means!



Classification with high-dimensional data

- Linear / Quadratic discriminant analysis
- Logistic regression
- Support Vector Machines



$$C = \{0, 1\}$$

$$P(C=1|x) = \beta_0 + \beta^T x$$

$\hat{C} \quad \hat{\beta}$

$(0, 1) \quad (-\infty, \infty)$

$$\log\left(\frac{p}{1-p}\right)$$

\hat{m}

$(-\infty, \infty)$

Logistic regression

- Assume we have two classes $C = \{0, 1\}$
- $Y \sim \text{Bernoulli}(p)$, and we would like to model p
- What is wrong with assuming a linear model for p , i.e.
 $p = \beta_0 + \beta^T x$

- In logistic regression we model: $\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta^T x$

- We need to estimate β_0 and β (How?)

- Then, we can solve for p , namely

$$p = \frac{\exp(\beta_0 + \beta^T x)}{1 + \exp(\beta_0 + \beta^T x)}$$

- Choose $C = 1$ if $p > 0.5$

$$\binom{n}{h} p^h (1-p)^{n-h}$$

$h \log(p) + (n-h) \log(1-p)$

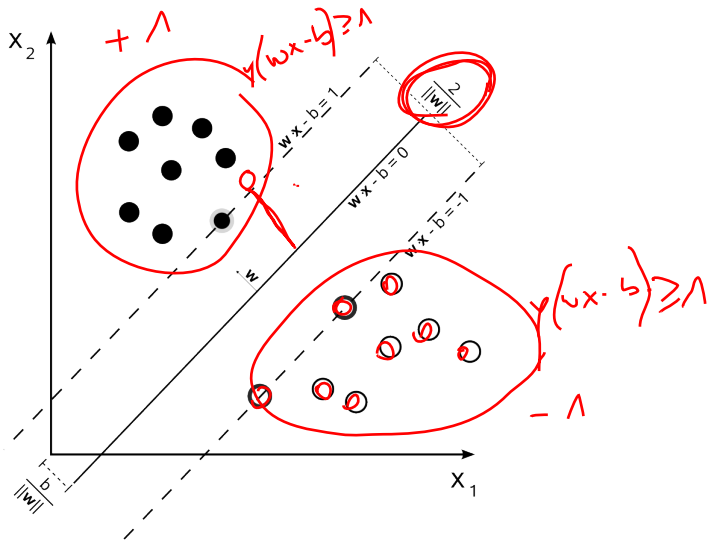


Figure from https://en.wikipedia.org/wiki/Support_vector_machine

- Given training data: $(x_1, y_1), \dots, (x_n, y_n)$, with $x_i \in \mathbb{R}^p$, $y_i \in \{-1, 1\}$
- If perfect classification is possible, determine hyperplane ($wx - b = 0$) that maximizes distance to the nearest point x_i from each group:

$$\text{minimize } \|w\|_2 \quad \text{such that } \underline{y_i(wx_i - b) \geq 1} \quad \text{for all } i.$$

- Given training data: $(x_1, y_1), \dots, (x_n, y_n)$, with $x_i \in \mathbb{R}^p$, $y_i \in \{-1, 1\}$
- If perfect classification is possible, determine hyperplane ($w x - b = 0$) that maximizes distance to the nearest point x_i from each group:

$$\text{minimize } \|w\|_2 \quad \text{such that } y_i(w x_i - b) \geq 1 \quad \text{for all } i.$$

- If perfect classification is not possible, determine hyperplane ($w x - b = 0$) that maximizes distance to the nearest point x_i from each group and minimizes sum of classification errors φ_i :

$$\text{minimize } \|w\|_2 + \lambda \sum_{i=1}^n \varphi_i \quad \text{such that } \underline{y_i(w x_i - b) \geq 1 - \varphi_i} \quad \text{for all } i.$$

Summary: Classification

- **Bayesian decision theory**

- Know probability distribution of the categories
- Do not even need training data
- Can design optimal classifier

- **Linear / Quadratic Discriminant Analysis (LDA / QDA)**

- Shape of probability distribution of categories is known (Gaussian)
- Need to estimate parameters of probability distribution
- Need training data

← logistic regression

- **Support Vector Machine (SVM)**

- Non-parametric method, i.e., no probability distribution
- Need to estimate parameters of discriminant function
- Labeled data, need training data

- **Clustering**

- Unlabeled data

Quality of classification

2 approaches:

- separate training data from test data
- cross validation, e.g. leave-one-out cross validation, where every sample is the test case once, the rest is the training data

Measures for prediction error:

- Build a **confusion matrix**

	Truth = 0	Truth = 1	Truth = 2
Estimate = 0	23	7	6
Estimate = 1	3	27	4
Estimate = 2	3	1	26

- Error rate = $1 - \text{sum}(\text{diagonal entries}) / (\text{total number of samples})$

Chapter 4 in

- T. Hastie, R. Tibshirani, & J. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2009.

I. Goodfellow, Y. Bengio, & A. Courville. *Deep Learning*. MIT Press, 2016.

M. Belkin, D. Hsu, S. Ma, & S. Mandal. Reconciling modern machine-learning practice and the classical bias-variance trade-off. *PNAS* 116, 2019.