

CS 412

FEB 6TH – LOGISTIC REGRESSION / SVM

HTF – CHAPTER 12

HW1 Review

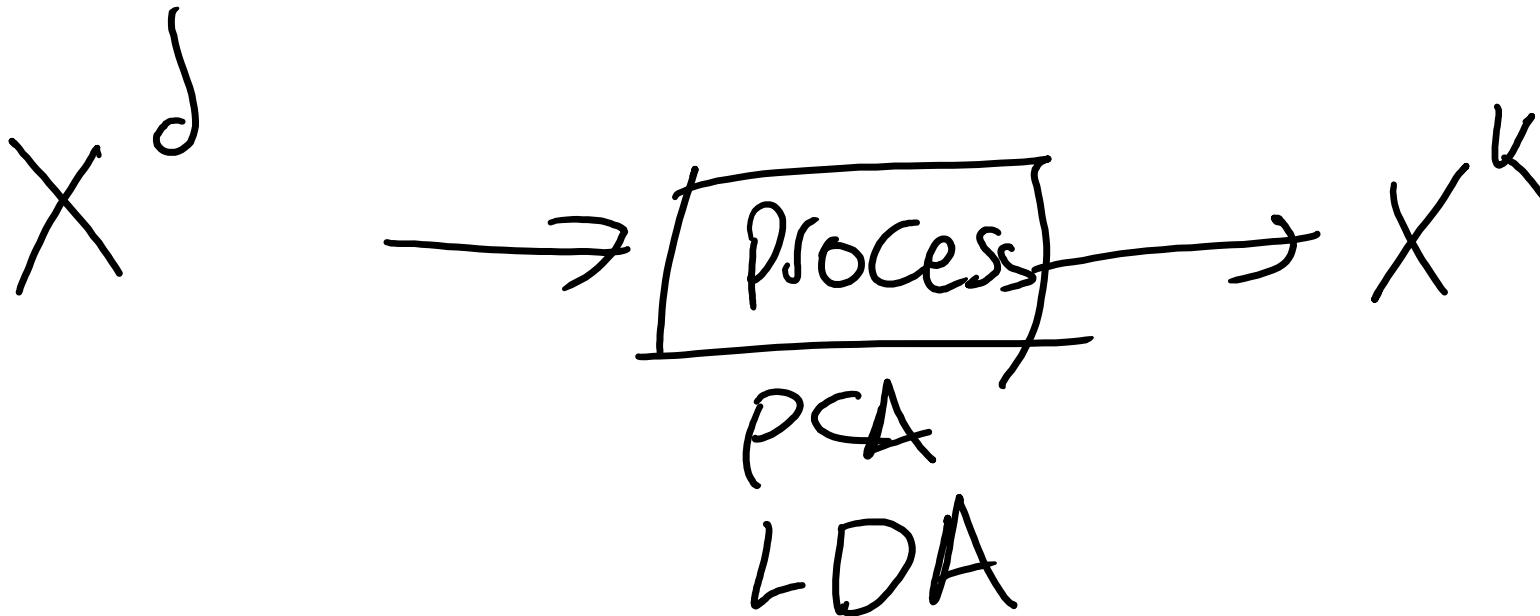
Important takeaways from HW1

- `cross_val_score()`
- Unusual graphs?
- What values of k seemed best?
- Where there any mistakes in the reasoning of the process?

Feature Extraction

Rather than selecting a subset of features, we can procedurally generate new features

- What dimensions in the data are best at explaining?
- They'll often be combinations of features, and as usual, we want to offload the selection process as much as possible.



Principal Components Analysis

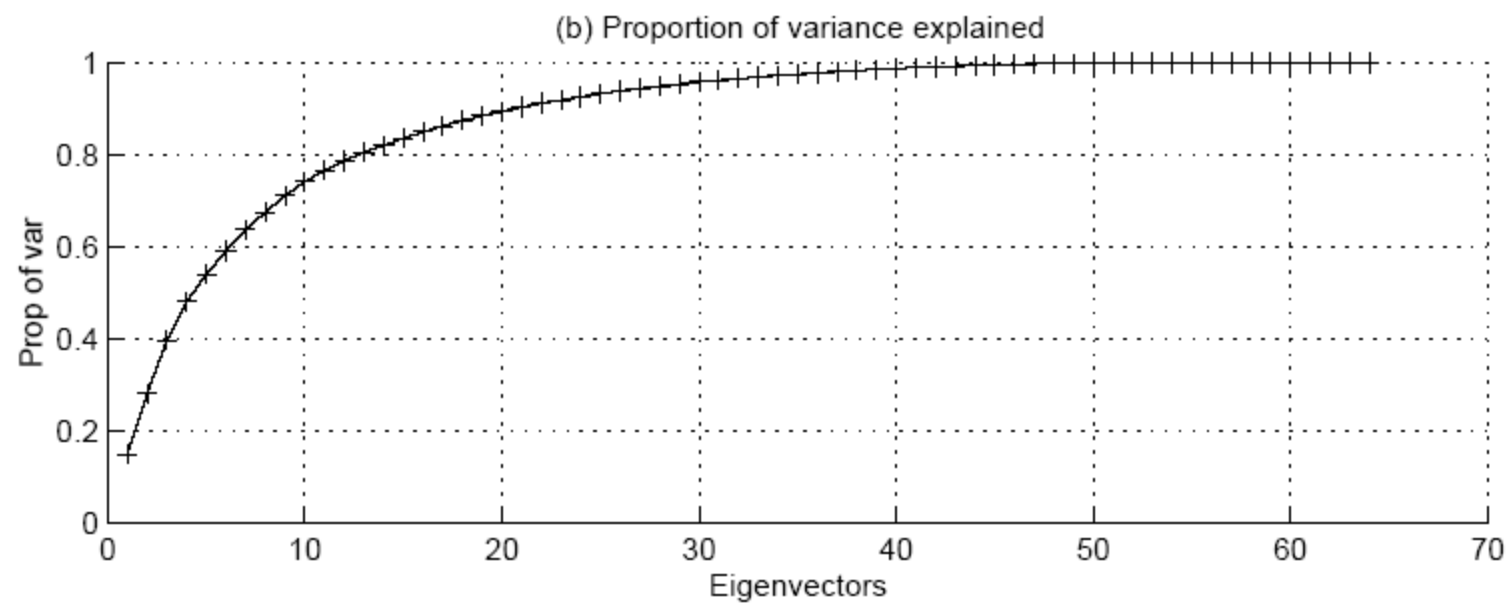
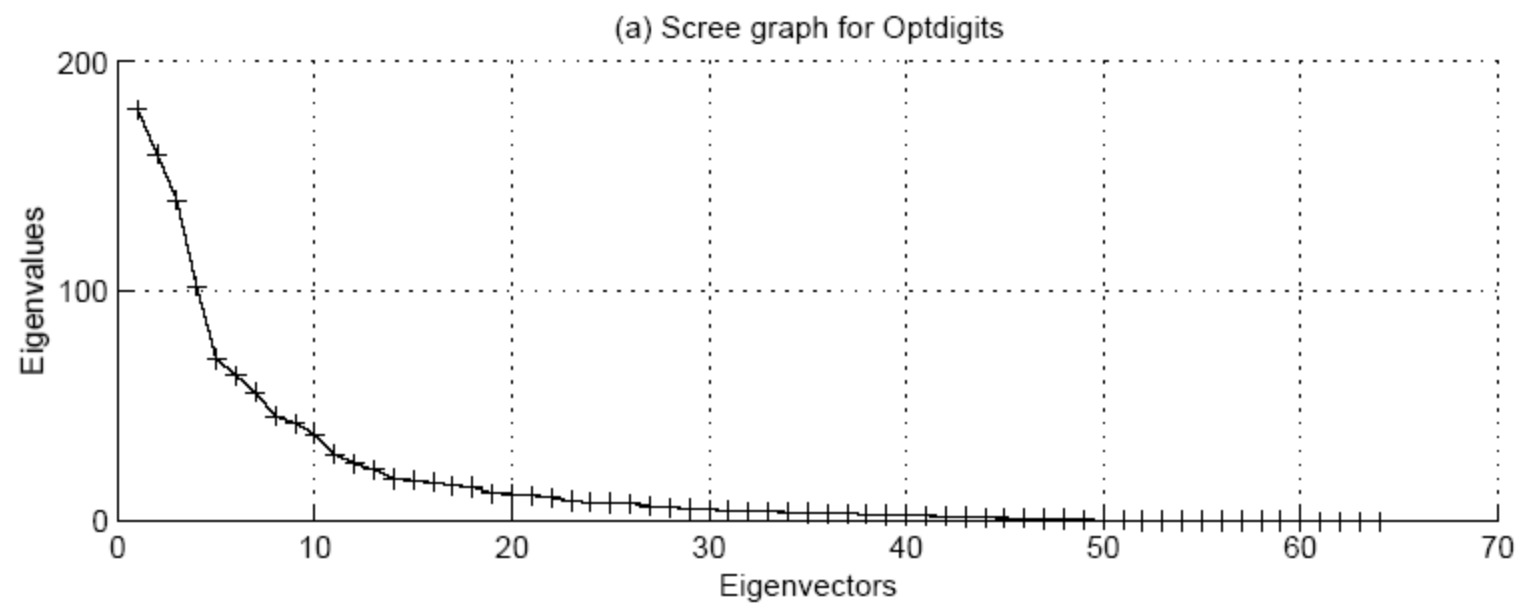
Find a low-dimensional space such that when \mathbf{x} is projected there, information loss is minimized.

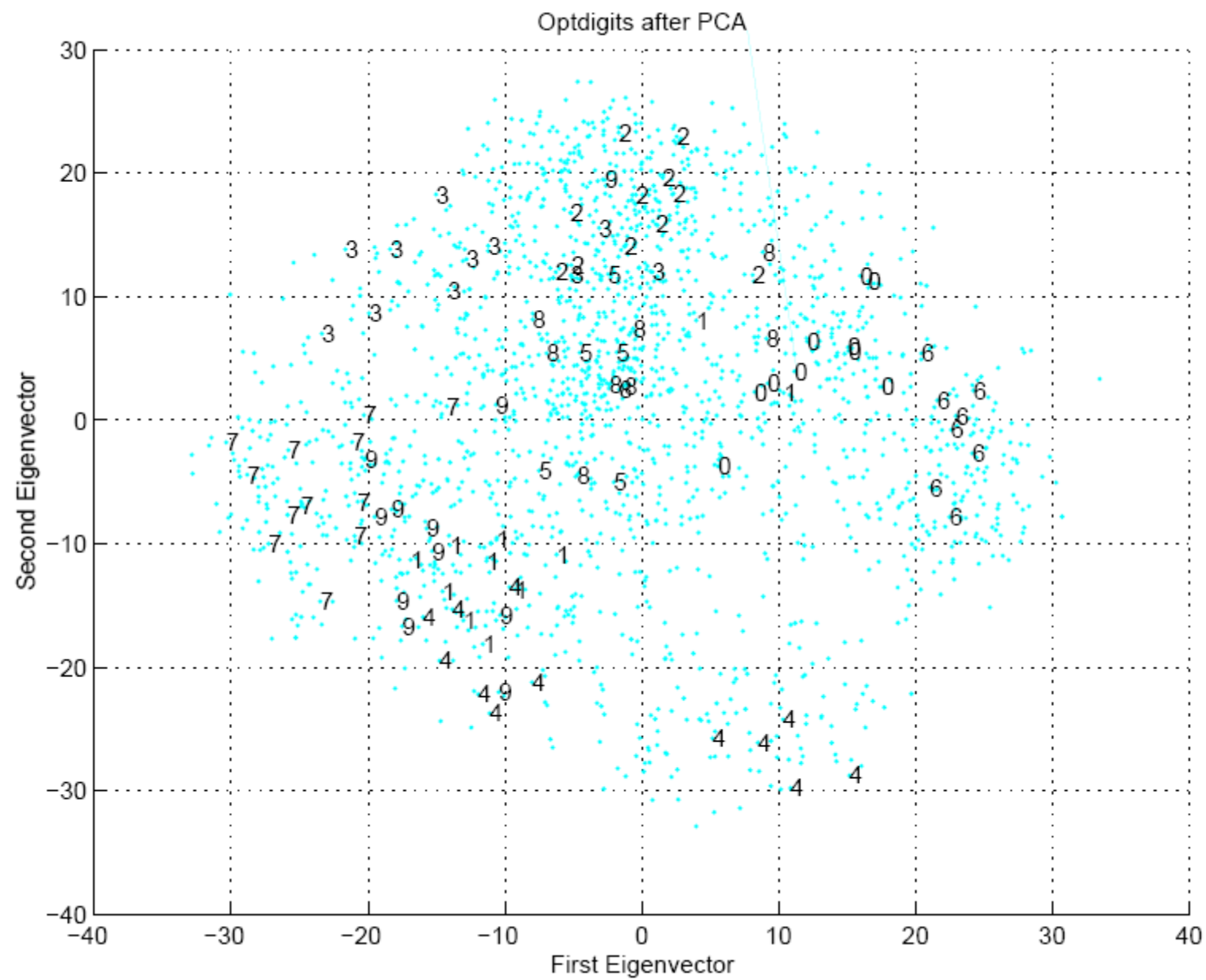
The projection of \mathbf{x} on the direction of \mathbf{w} is: $z = \mathbf{w}^T \mathbf{x}$

Find \mathbf{w} such that $\text{Var}(z)$ is maximized

$$\begin{aligned}\text{Var}(z) &= \text{Var}(\mathbf{w}^T \mathbf{x}) = \text{E}[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})^2] \\ &= \text{E}[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})] \\ &= \text{E}[\mathbf{w}^T (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{w}] \\ &= \mathbf{w}^T \text{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] \mathbf{w} = \mathbf{w}^T \Sigma \mathbf{w}\end{aligned}$$

where $\text{Var}(\mathbf{x}) = \text{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] = \Sigma$





Linear Discriminant Analysis

Find a low-dimensional space such that when x is projected, classes are well-separated.

Find w that maximizes the separation

Linear Discriminant Analysis

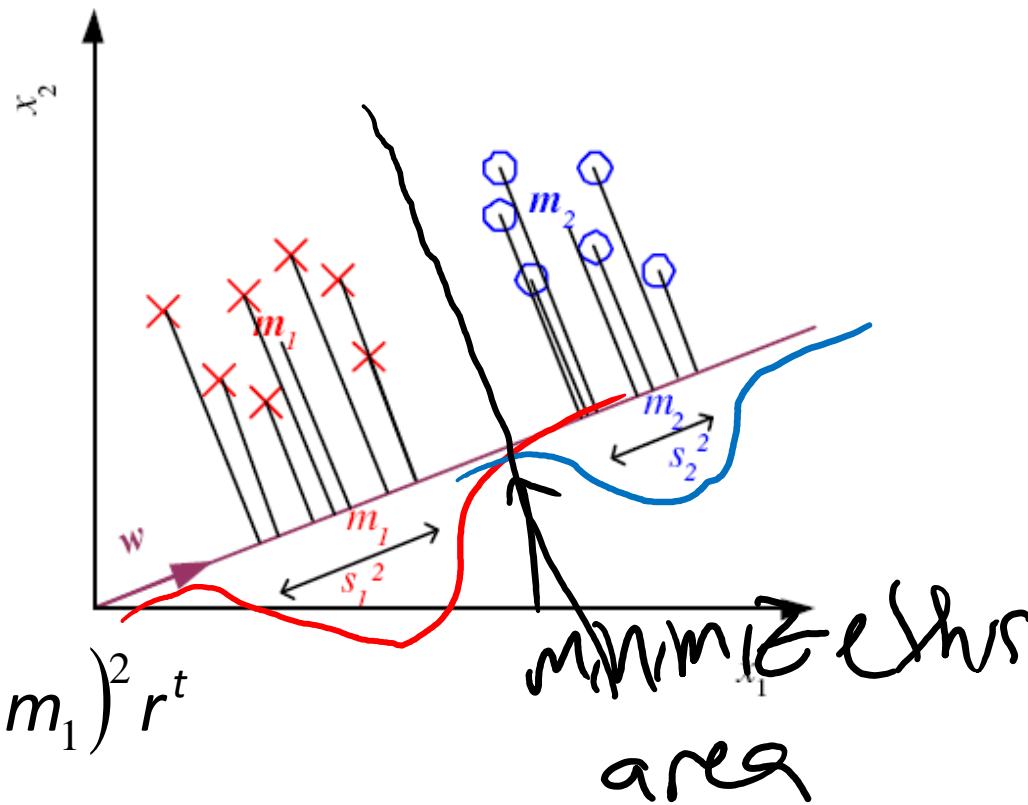
Find a low-dimensional space such that when is projected, classes are well-separated.

Find w that maximizes the separation

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$$

$$m_1 = \frac{\sum_t \mathbf{w}^T \mathbf{x}^t r^t}{\sum_t r^t}$$

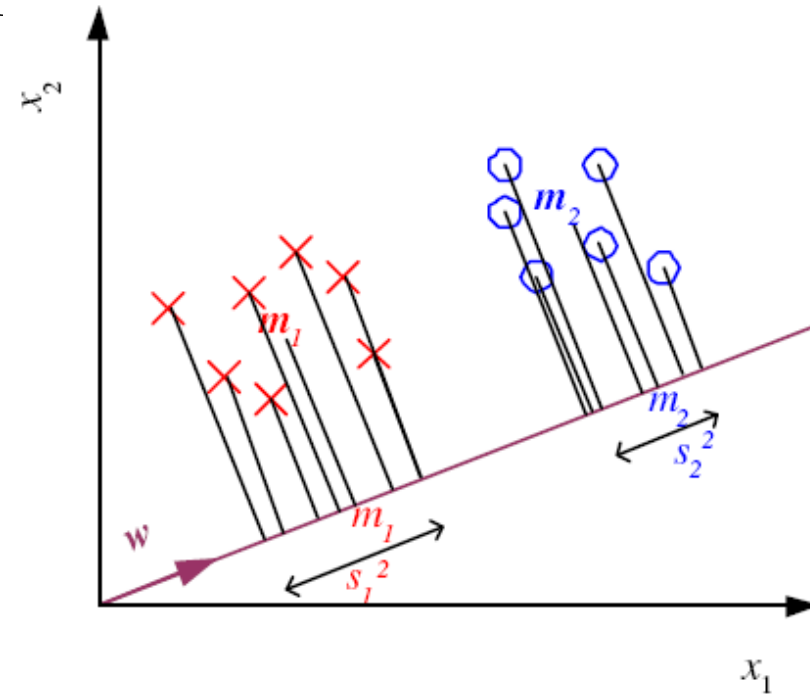
$$s_1^2 = \sum_t (\mathbf{w}^T \mathbf{x}^t - m_1)^2 r^t$$



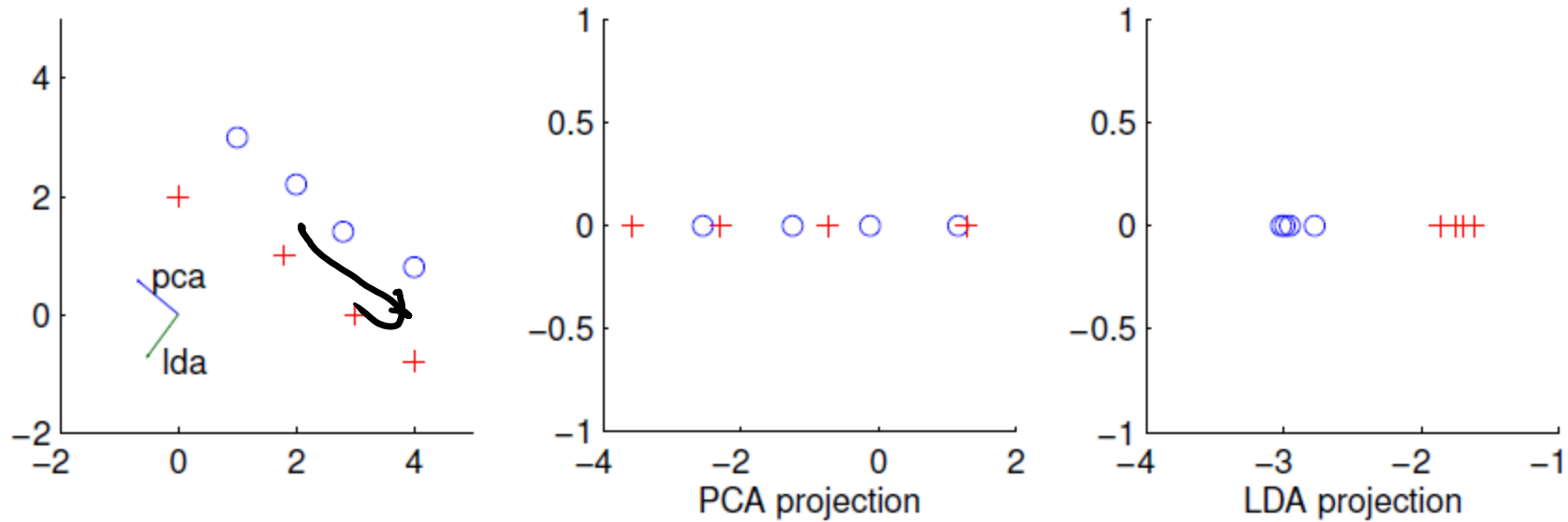
Linear Discriminant Analysis

Find a low-dimensional space such that when is projected, classes are well-separated.

Find w that maximizes the separation



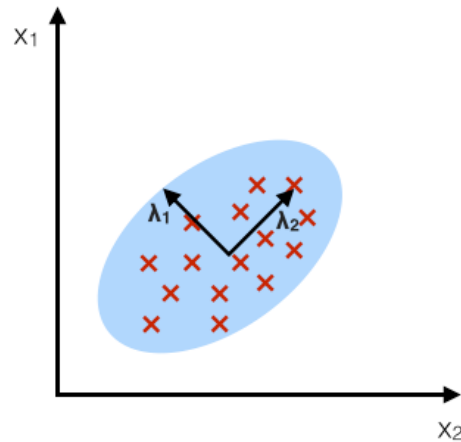
PCA vs LDA



PCA vs LDA

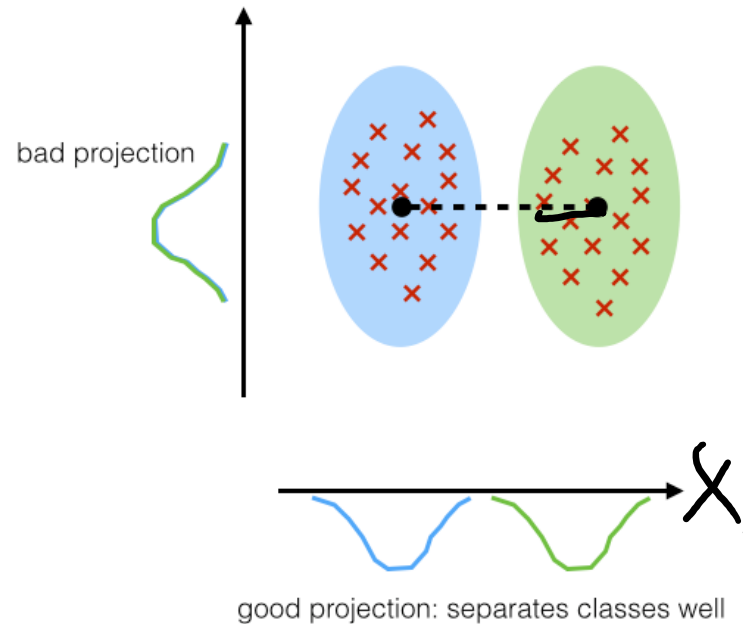
PCA:

component axes that maximize the variance



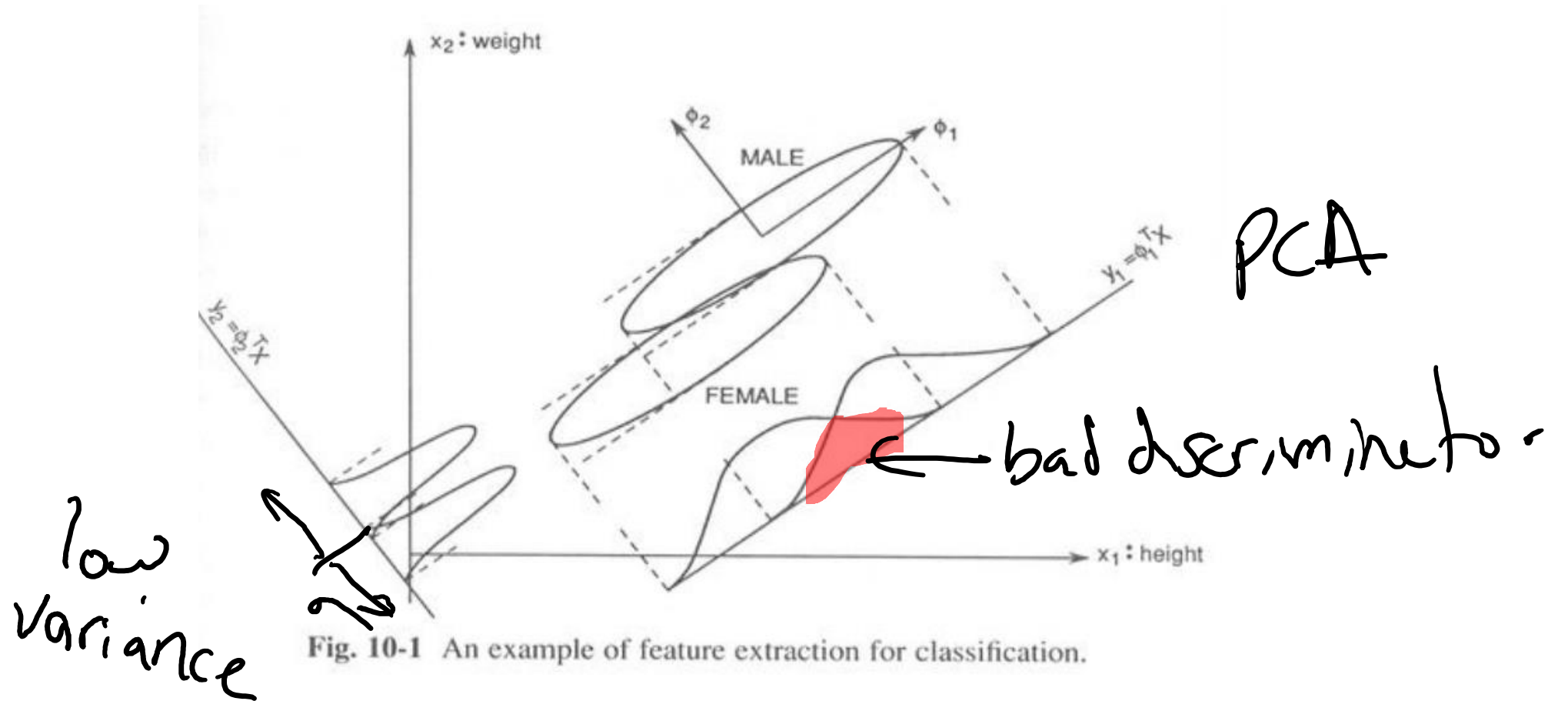
LDA:

maximizing the component axes for class-separation

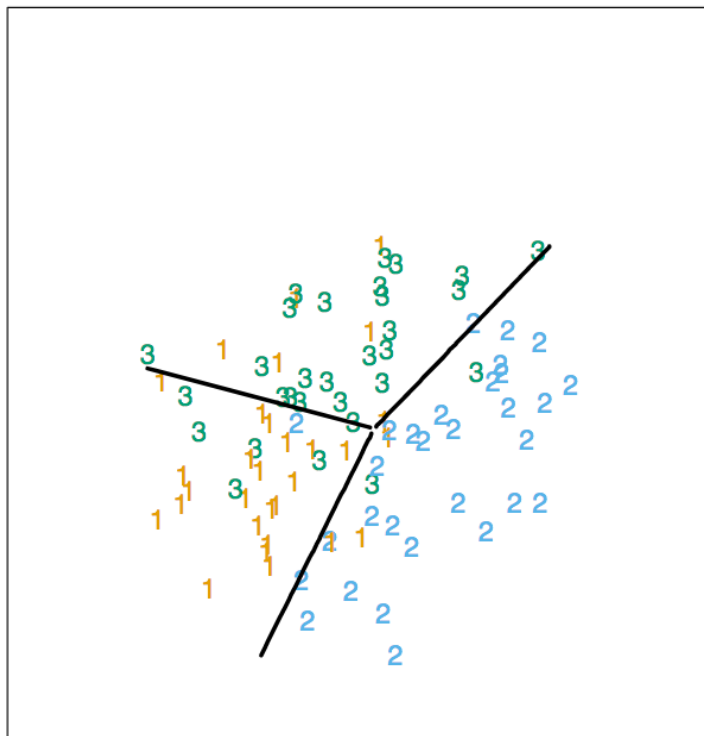
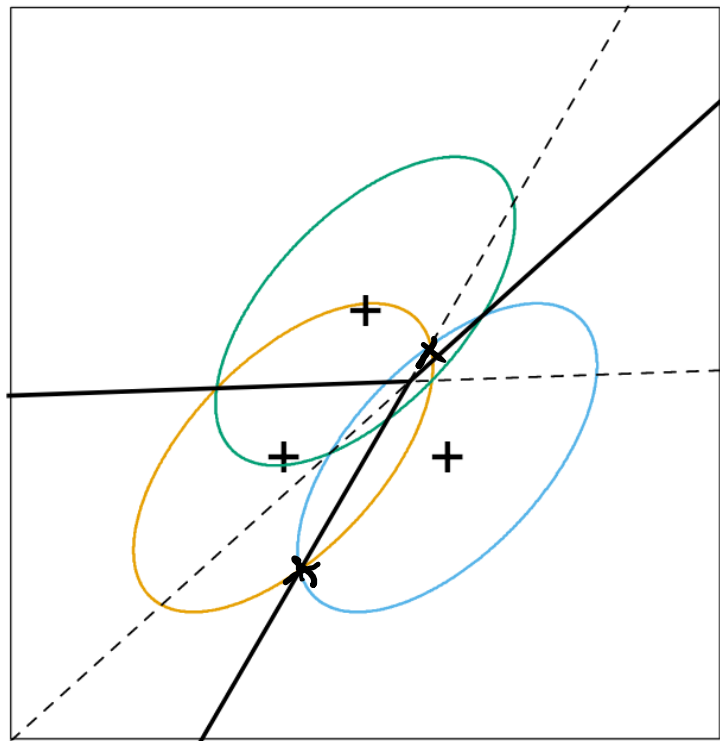


PCA and LDA are not limited to the axes

PCA vs LDA



identically distributed



Discriminate analysis

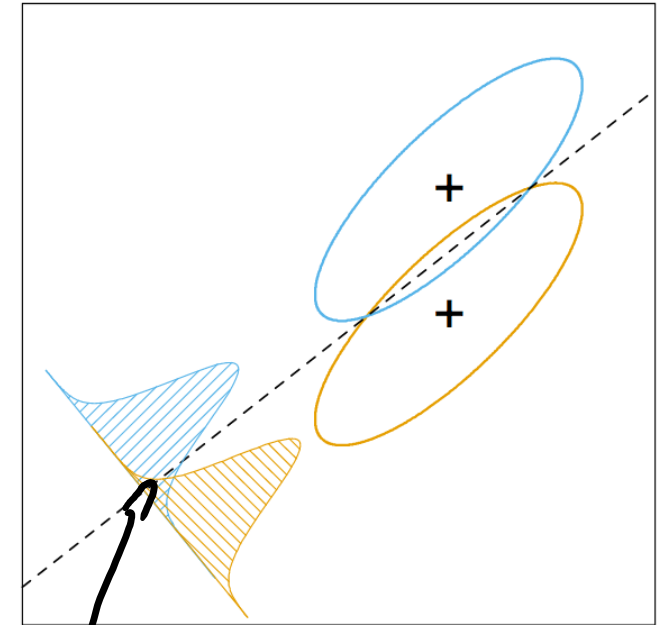
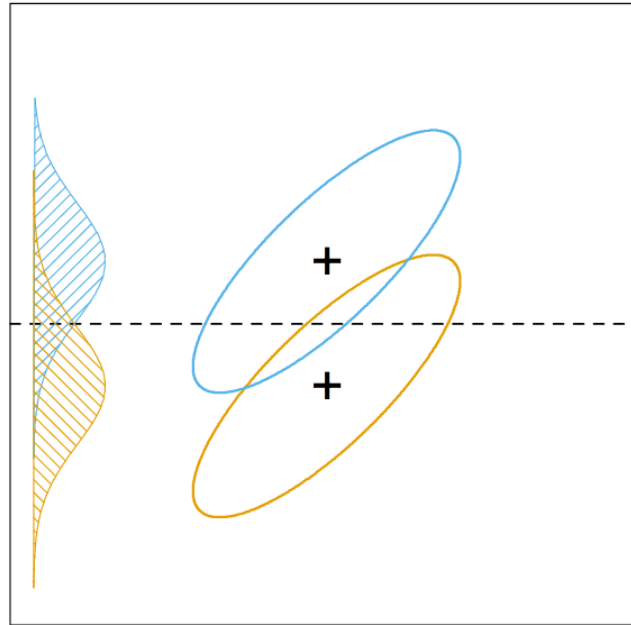
Suppose that the data is randomly distributed in \mathbb{R}^p relative to some multivariate Gaussian distribution.

- Find the center of each distribution region and then separate the regions based on what minimizes the probability of overlap
- Assume that each of the regions shares the same distribution

Discriminate analysis

To minimize overlap, we want to discriminate between the classes along the axis that best discriminates the data

Unlike linear regression we want to divide the region into relevant spaces

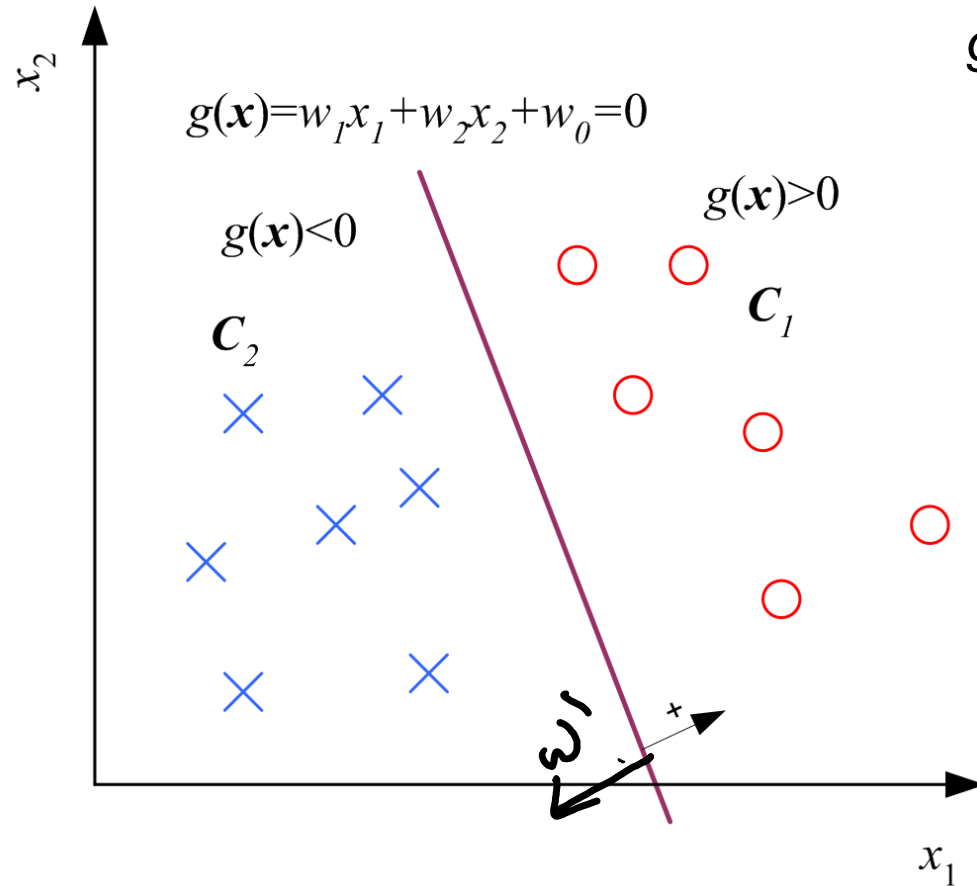


minimal overlap

Two Classes

Polynomial kernel $\phi(x)$

$x_1, x_2, x_1^2, x_2^2, x_1x_2$



choose $\begin{cases} C_1 & \text{if } g(\mathbf{x}) > 0 \\ C_2 & \text{otherwise} \end{cases}$

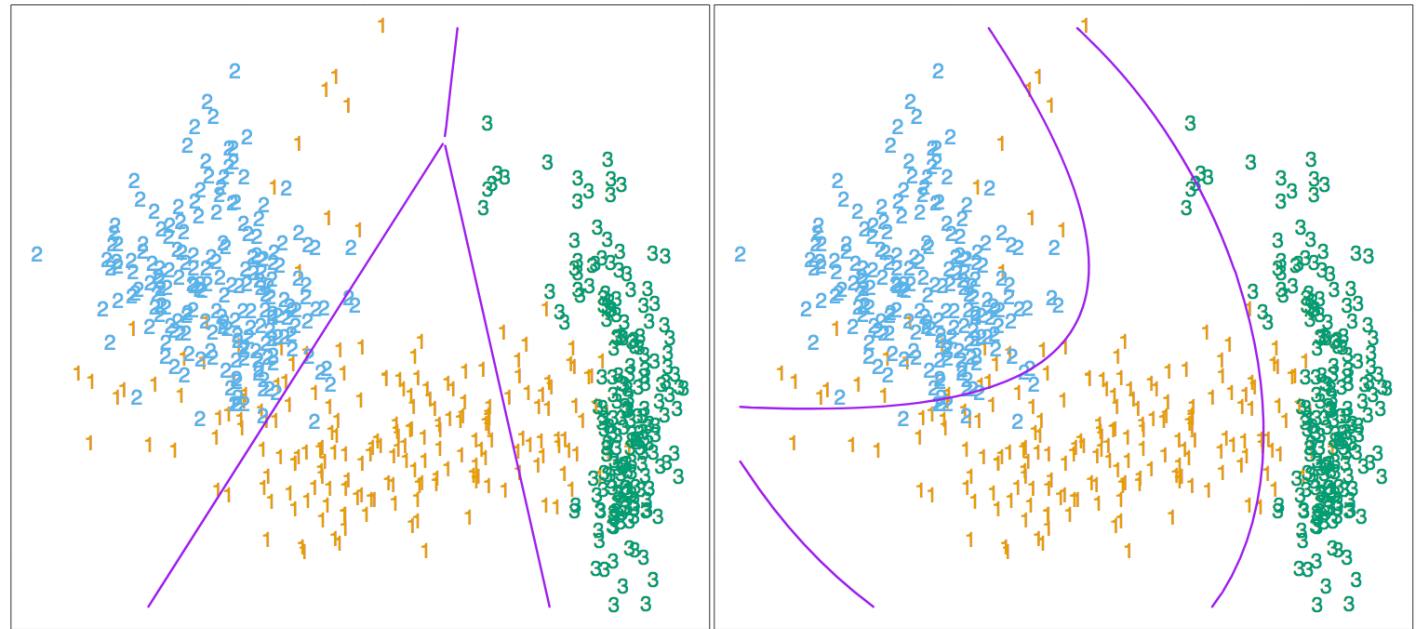
$C_2 \quad g'(\mathbf{x}) > 0$
 $C_1 \quad \text{otherwise}$

Quadratic discriminants

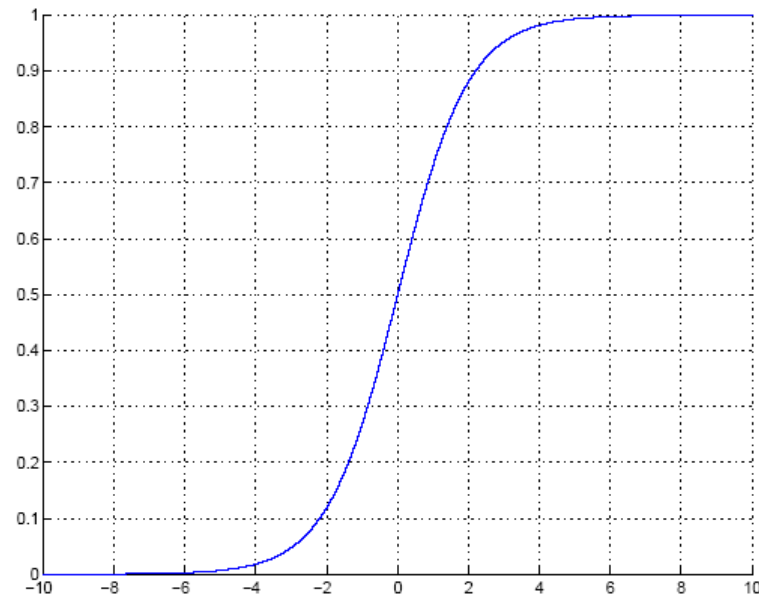
Here, we can see two different models

- X_1, X_2
- $X_1, X_2, X_1X_2, X_1^2, X_2^2$

For only two variables, it is easy to generate the interaction/polynomial terms



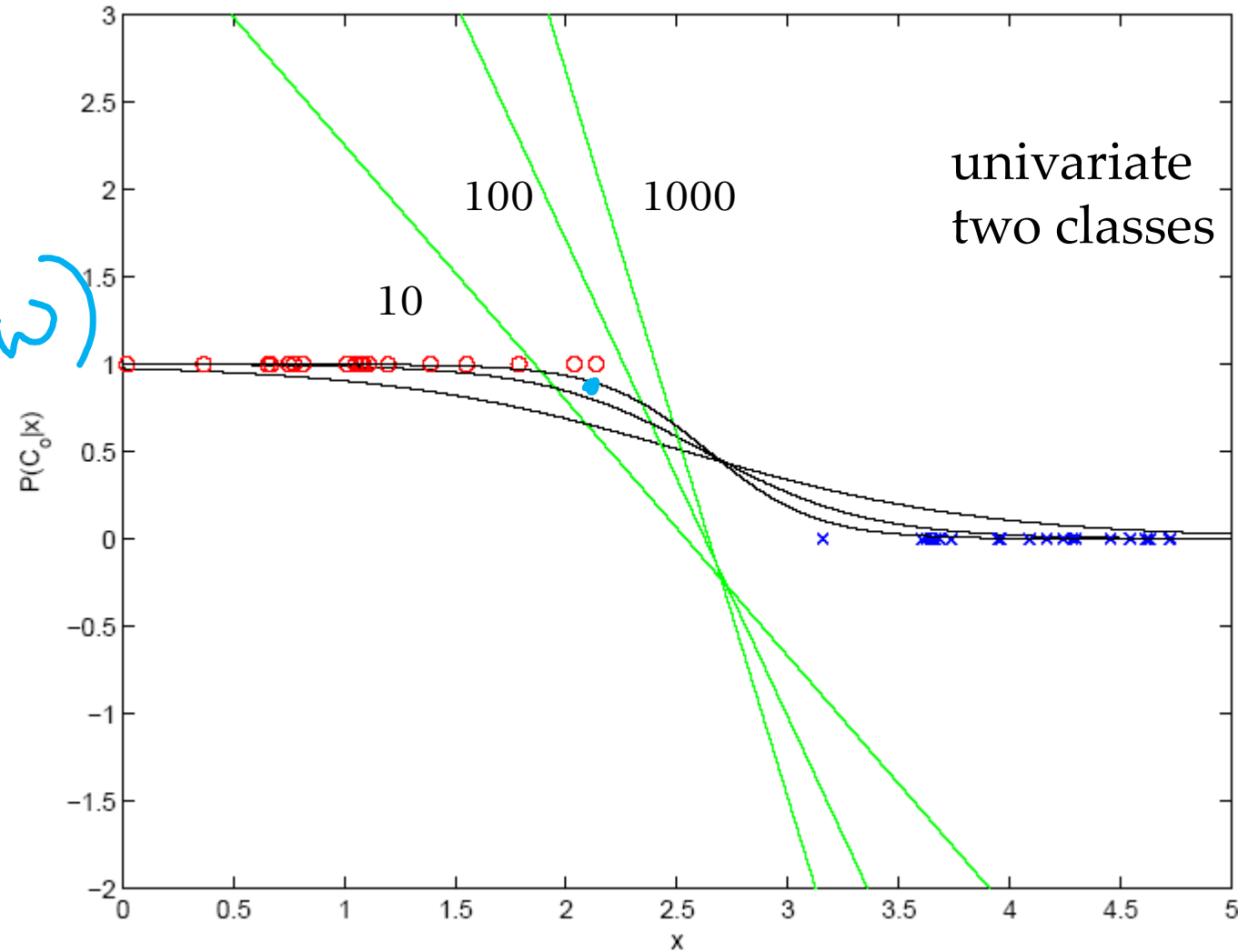
Sigmoid (Logistic) Function



Calculate $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ and choose C_1 if $g(\mathbf{x}) > 0$, or
Calculate $y = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + w_0)$ and choose C_1 if $y > 0.5$

$$= \text{sigmoid}(a), \text{ where } a = \mathbf{w}^T \mathbf{x} + w_0 \quad \frac{dy}{da} = y(1-y)$$

after 10, 100, 1000 iterations



$\text{Sigmoid}(xw)$
11
0.9

$g(x) = \text{Sigmoid}(x^T w)$
 $g(x) > 0.5$
red '0'

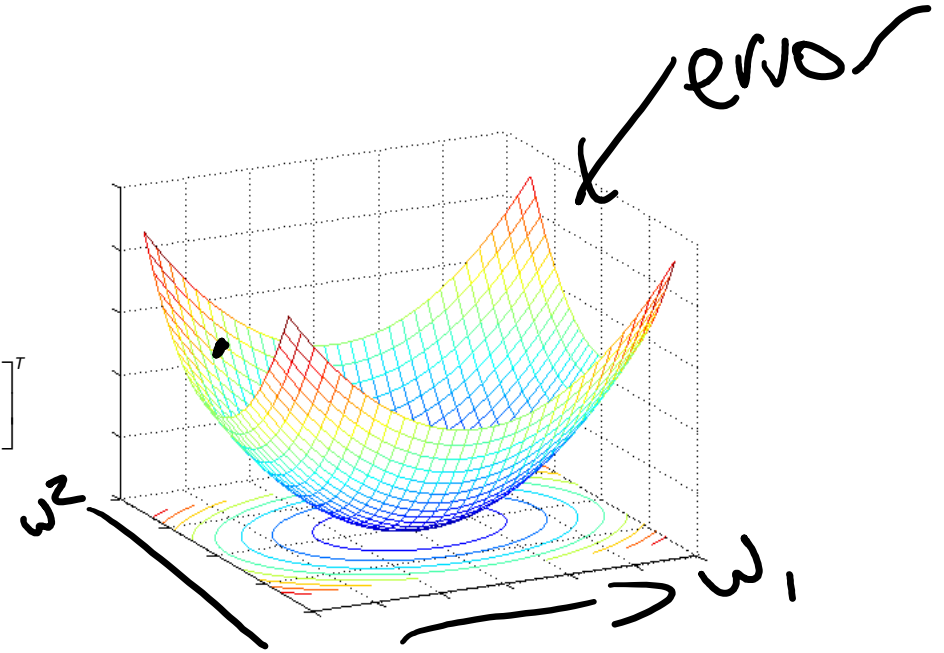
Gradient-Descent

$E(w | X)$ is error with parameters w on sample X

$$w^* = \arg \min_w E(w | X)$$

Gradient

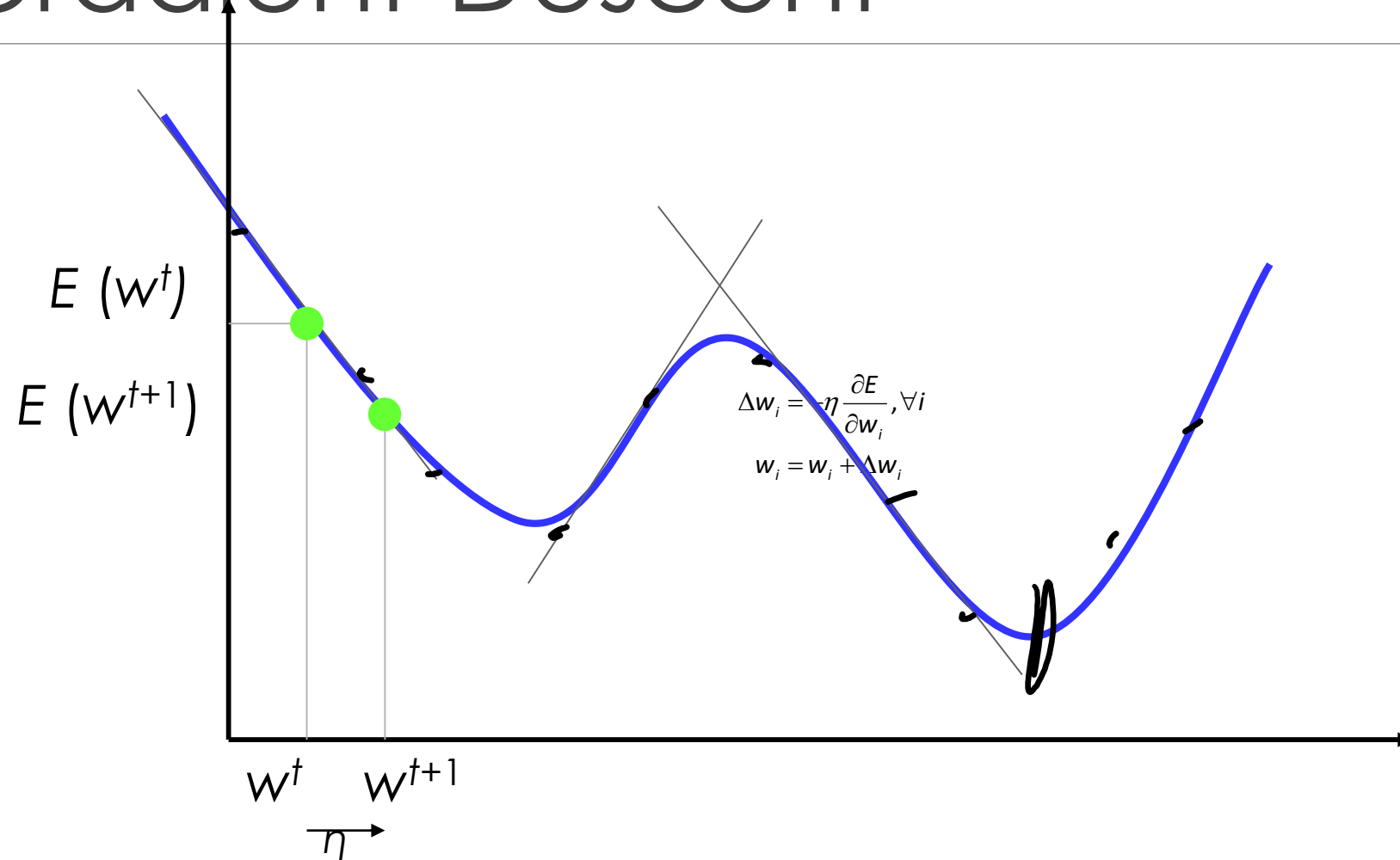
$$\nabla_w E = \left[\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_d} \right]^T$$



Gradient-descent:

Starts from random w and updates w iteratively in the negative direction of gradient

Gradient-Descent

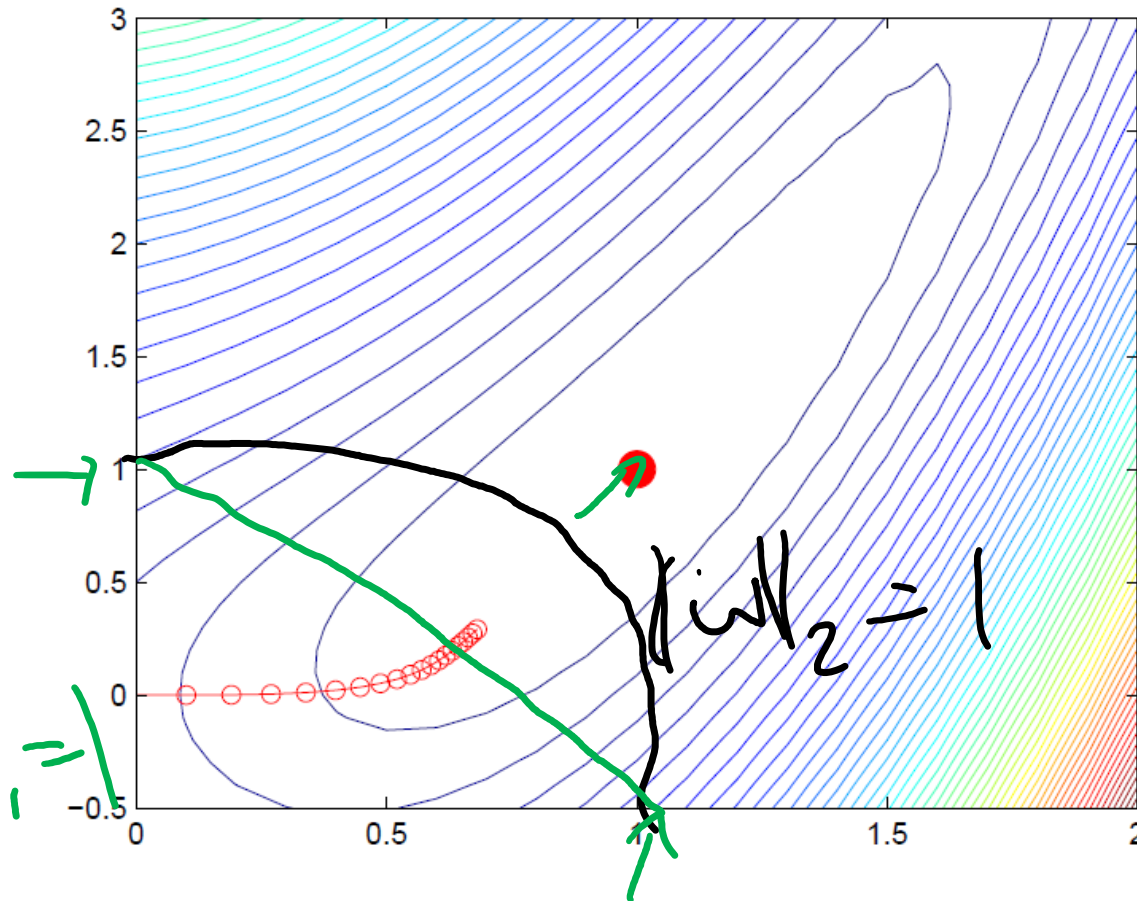


Learning for logistic regression

$$\eta = 0.1$$

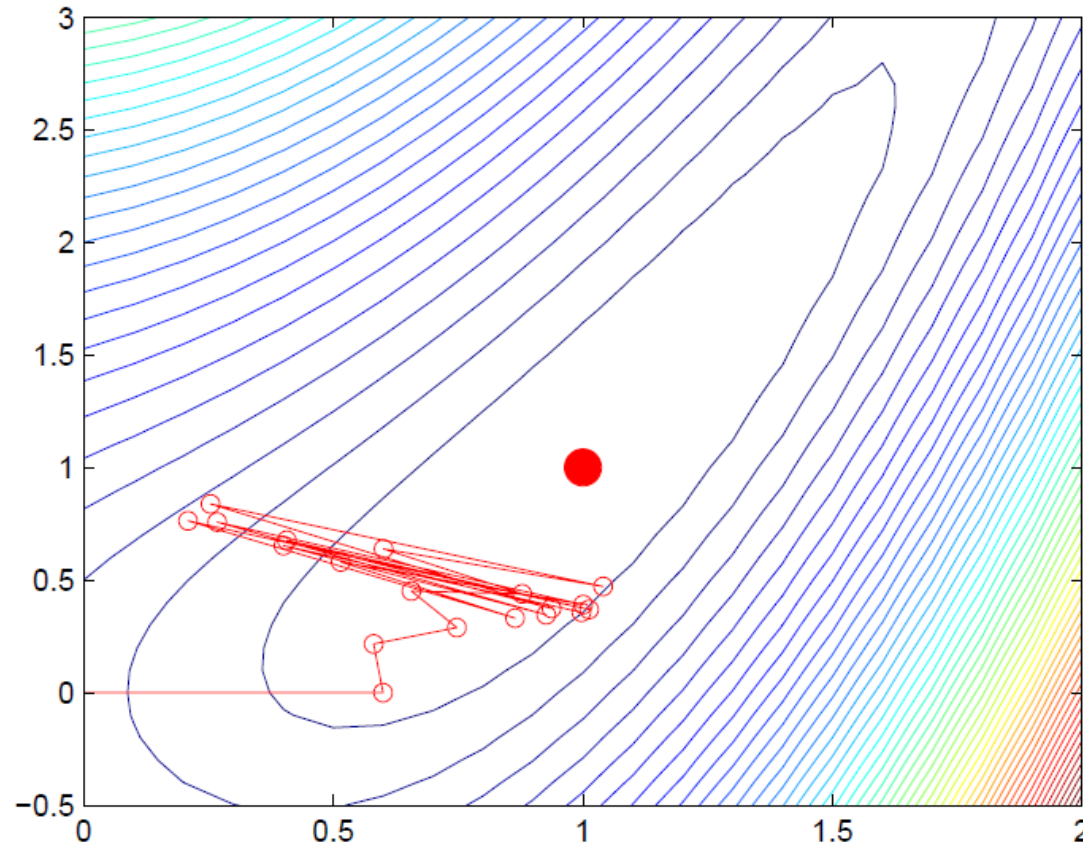
η = learning rate
(e/a)

$$\|w\|_1 = 1$$



Learning for logistic regression

$$\eta = 0.6$$



Logistic regression and overfitting

Overfitting

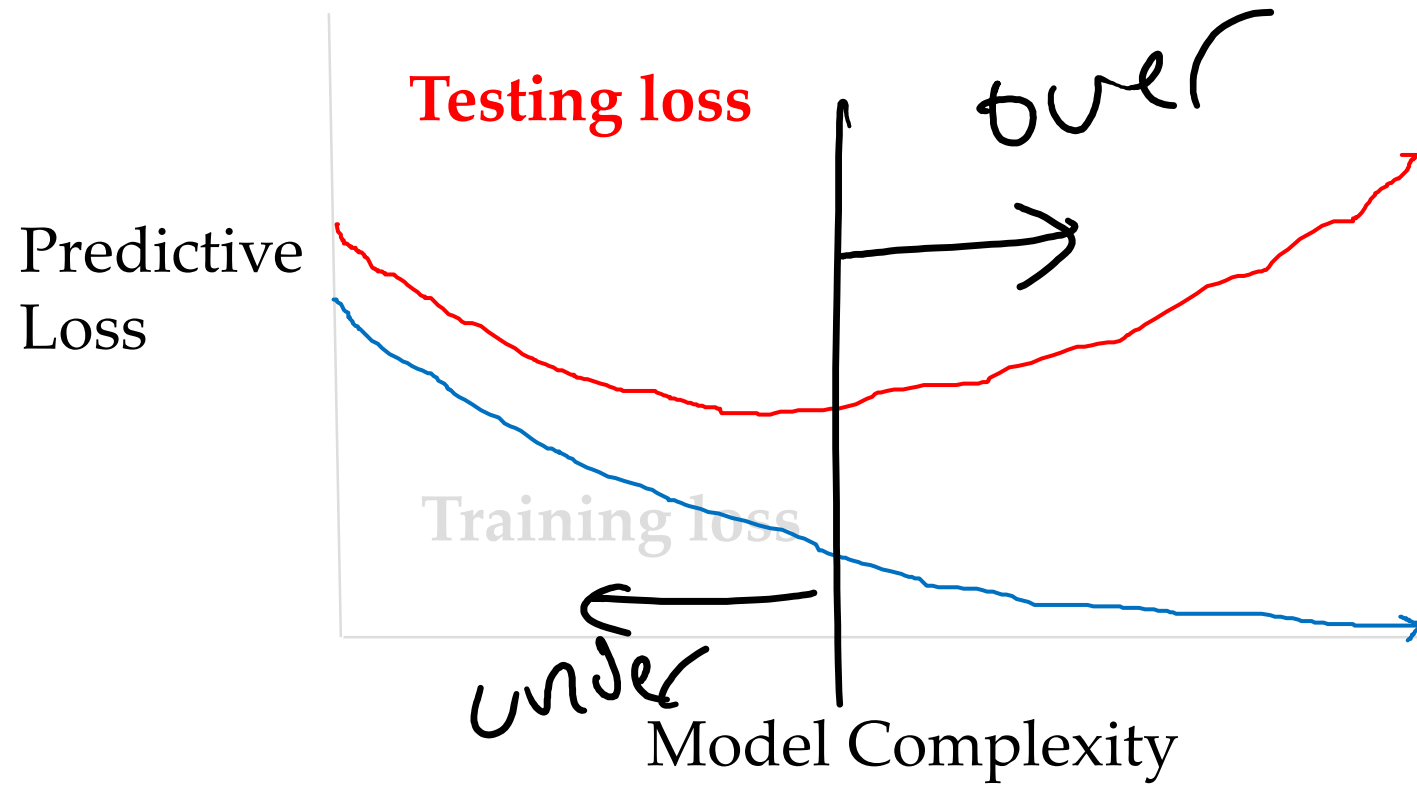
⇒ data points

- Occurs when very few instances and feature space is high dimensional

To avoid, a common approach is defining a prior on w

- Corresponds to Regularization
- Helps with avoiding large weights
- “Pushes” parameters to zero

Overfitting



Need to prevent complex hypotheses

Overfitting

- Occurs when very few instances and feature space is high dimensional

Idea #1: Restrict the number of features considered

- Cross-validation

Idea #2: Penalize complex hypotheses in the model search

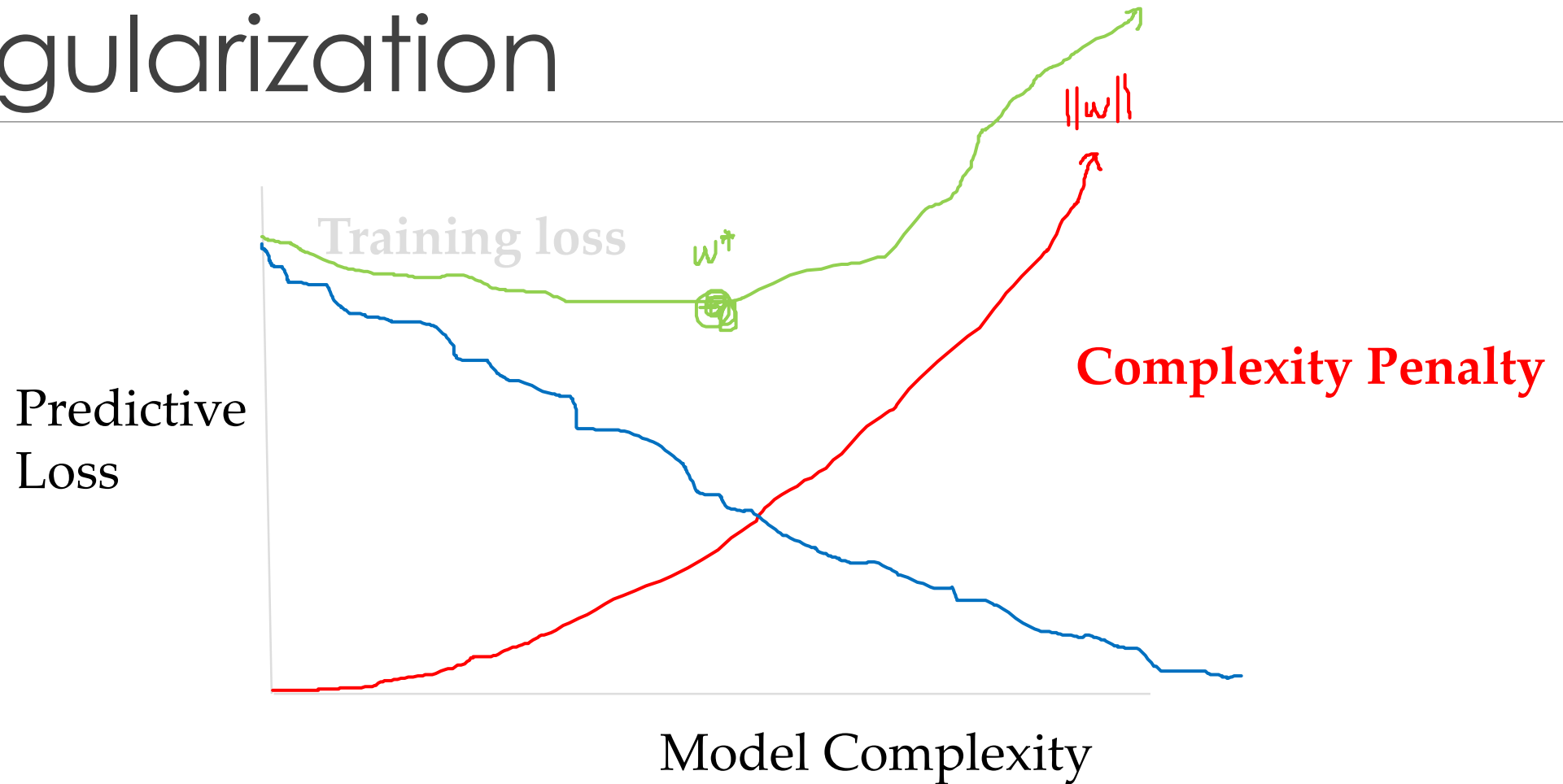
- Regularization!

$$g(x) = 1000x_1 - 1500x_2 + 100$$

$$h(x) = 10x_1 - 15x_2 + 1$$

$$= 0.001x_1 - 0.0015x_2 + 0.0001$$

Regularization



Regularization

Recall the objective of logistic regression:

$$E(\mathbf{w}, w_0 | \mathcal{X}) = - \sum_t r^t \log y^t + (1 - r^t) \log (1 - y^t)$$

Sigmoid function

L2 regularization

$$\underset{\mathbf{w}}{\operatorname{argmin}} E(\mathbf{w}, w_0 | X) + \lambda \sum_i w_i^2$$

L1 regularization

$$\underset{\mathbf{w}}{\operatorname{argmin}} E(\mathbf{w}, w_0 | X) + \lambda \sum_i |w_i|$$

λ = regularization parameter

$\lambda > 0$ is a weight, chosen by, e.g., cross validation