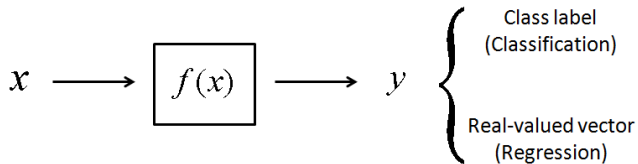


# Machine Learning Basics

September 11-12, 2018

# Machine Learning



Object recognition

{dog, cat, horse, flower, ...}



Low-resolution image

Super resolution



High-resolution image

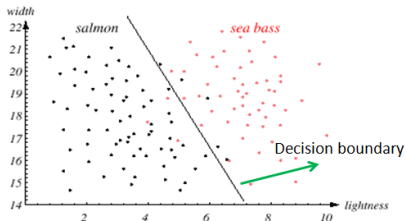
# Classification

- $f(\mathbf{x})$  predicts the category that  $\mathbf{x}$  belongs to

$$f: \mathcal{R}^D \rightarrow \{1, \dots, K\}$$

- $f(\mathbf{x})$  is decided by the decision boundary
- As an variant,  $f$  can also predict the probability distribution over classes given  $\mathbf{x}$ ,  $f(\mathbf{x}) = P(y|\mathbf{x})$ . The category is predicted as

$$y^* = \arg \max_k P(y = k|\mathbf{x})$$



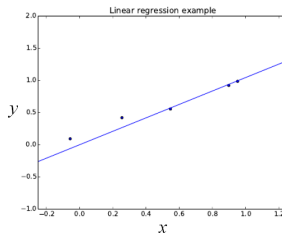
# Regression

- Predict real-valued output

$$f : \mathcal{R}^D \rightarrow \mathcal{R}^M$$

- Example: linear regression

$$y = \mathbf{w}^t \mathbf{x} = \sum_{d=1}^D w_d x_d + w_0$$



# Training

- Training: estimate the parameters of  $f$  from  $\{(\mathbf{x}_i^{(\text{train})}, y_i^{(\text{train})})\}$ 
  - Decision boundary, parameters of  $P(y|\mathbf{x})$ , and  $\mathbf{w}$  in linear regression
- Optimize an objective function on the training set. It is a performance measure on the training set and could be different from that on the test set.
  - Mean squared error (MSE) for linear regression

$$\text{MSE}_{\text{train}} = \frac{1}{N} \sum_i ||\mathbf{w}^t \mathbf{x}_i^{(\text{train})} - y_i^{(\text{train})}||_2^2$$

- Cross entropy (CE) for classification

$$\text{CE}_{\text{train}} = \frac{1}{N} \sum_i \log P(y = y_i^{(\text{train})} | \mathbf{x}_i^{(\text{train})})$$

Why not use classification errors  $\#\{f(\mathbf{x}_i^{(\text{train})}) \neq y_i^{(\text{train})}\}$ ?

- The choice of the objective function should be good for optimization
- Take linear regression as an example

$$\nabla_{\mathbf{w}} \text{MSE}_{\text{train}} = 0$$

$$\Rightarrow \nabla_{\mathbf{w}} ||\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}||_2^2 = 0$$

$$\mathbf{w} = (\mathbf{X}^{(\text{train})t} \mathbf{X}^{(\text{train})})^{-1} \mathbf{X}^{(\text{train})t} \mathbf{y}^{(\text{train})}$$

where  $\mathbf{X}^{(\text{train})} = [\mathbf{x}_1^{(\text{train})}, \dots, \mathbf{x}_N^{(\text{train})}]$  and  $\mathbf{y}^{(\text{train})} = [y_1^{(\text{train})}, \dots, y_N^{(\text{train})}]$ .

# Generalization

- We care more about the performance of the model on new, previously unseen examples
- The training examples usually cannot cover all the possible input configurations, so the learner has to generalize from the training examples to new cases
- Generalization error: the expected error over **ALL** examples
- To obtain theoretical guarantees about generalization of a machine learning algorithm, we assume all the samples are drawn from a distribution  $p(\mathbf{x}, y)$ , and calculate generalization error (GE) of a prediction function  $f$  by taking expectation over  $p(\mathbf{x}, y)$

$$GE_f = \int_{\mathbf{x}, y} p(\mathbf{x}, y) \mathbf{Error}(f(\mathbf{x}), y)$$

# Generalization

- However, in practice,  $p(\mathbf{x}, y)$  is unknown. We assess the generalization performance with a test set  $\{\mathbf{x}_i^{(\text{test})}, y_i^{(\text{test})}\}$

$$\text{Performance}_{\text{test}} = \frac{1}{M} \sum_{i=1}^M \text{Error}(f(\mathbf{x}_i^{(\text{test})}), y_i^{(\text{test})})$$

- We hope that both test examples and training examples are drawn from  $p(\mathbf{x}, y)$  of interest, although it is unknown



# Capacity

- The ability of the learner (or called model) to discover a function taken from a family of functions. Examples:

- Linear predictor

$$y = wx + b$$

- Quadratic predictor

$$y = w_2x^2 + w_1x + b$$

- Degree-10 polynomial predictor

$$y = b + \sum_{i=1}^{10} w_i x^i$$

- The latter family is richer, allowing to capture more complex functions
- Capacity can be measured by the number of training examples  $\{\mathbf{x}_i^{(\text{train})}, y_i^{(\text{train})}\}$  that the learner **could always fit**, no matter how to change the values of  $\mathbf{x}_i^{(\text{train})}$  and  $y_i^{(\text{train})}$

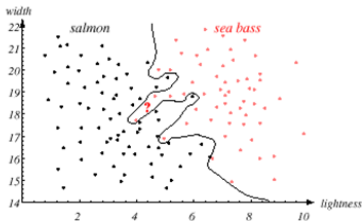
# Underfitting

- The learner cannot find a solution that fits training examples well
  - For example, use linear regression to fit training examples  $\{\mathbf{x}_i^{(\text{train})}, y_i^{(\text{train})}\}$  where  $y_i^{(\text{train})}$  is a quadratic function of  $\mathbf{x}_i^{(\text{train})}$
- Underfitting means the learner cannot capture some important aspects of the data
- Reasons for underfitting happening
  - Model is not rich enough
  - Difficult to find the global optimum of the objective function on the training set or easy to get stuck at local minimum
  - Limitation on the computation resources (not enough training iterations of an iterative optimization procedure)
- Underfitting commonly happens in deep learning with large scale training data and could be even a more serious problem than overfitting in some cases

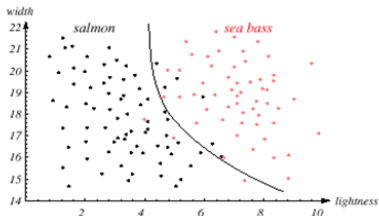
# Overfitting

- The learner fits the training data well, but loses the ability to generalize well, i.e. it has small training error but larger generalization error
- A learner with large capacity tends to overfit
  - The family of functions is too large (compared with the size of the training data) and it contains many functions which all fit the training data well.
  - Without sufficient data, the learner cannot distinguish which one is most appropriate and would make an arbitrary choice among these apparently good solutions
  - A separate validation set helps to choose a more appropriate one
  - In most cases, data is contaminated by noise. The learner with large capacity tends to describe random errors or noise instead of the underlying models of data (classes)

# Overfitting



Overly complex models lead to complicated decision boundaries. It leads to perfect classification on the training examples, but would lead to poor performance on new examples.



The decision boundary might represent the optimal tradeoff between performance on the training set and simplicity of classifier, therefore giving highest accuracy on new examples.

(Duda et al. Pattern Classification 2000)

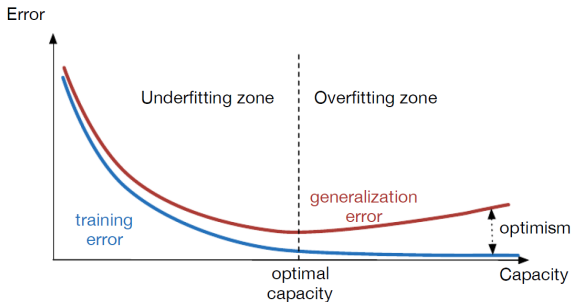
# Occam's Razor

- **The fundamental element of machine learning is the trade-off between capacity and generalization**
- Occam's Razor states that among competing functions that could explain the training data, one should choose the “simpler” one. Simplicity is the opposite of capacity.
- Occam's Razor suggests us pick the family of functions just large enough to leave only one choice that fits well the data.

# Optimal capacity

- Difference between training error and generalization error increases with the capacity of the learner
- Generalization error is a U-shaped function of capacity
- Optimal capacity capacity is associated with the transition from underfitting to overfitting
  - One can use a validation set to monitor generalization error empirically
- Optimal capacity should increase with the number of training examples

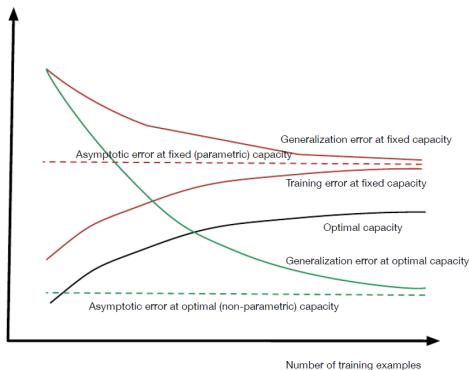
# Optimal capacity



Typical relationship between capacity and both training and generalization (or test) error. As capacity increases, training error can be reduced, but the optimism (difference between training and generalization error) increases. At some point, the increase in optimism is larger than the decrease in training error (typically when the training error is low and cannot go much lower), and we enter the overfitting regime, where capacity is too large, above the optimal capacity. Before reaching optimal capacity, we are in the underfitting regime.

(Bengio et al. Deep Learning 2014)

# Optimal capacity

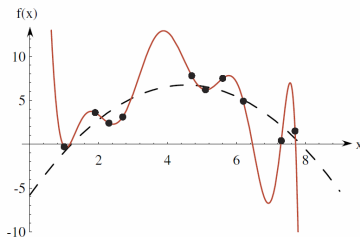


As the number of training examples increases, optimal capacity (bold black) increases (we can afford a bigger and more flexible model), and the associated generalization error (green bold) would decrease, eventually reaching the (non-parametric) asymptotic error (green dashed line). If capacity was fixed (parametric setting), increasing the number of training examples would also decrease generalization error (top red curve), but not as fast, and training error would slowly increase (bottom red curve), so that both would meet at an asymptotic value (dashed red line) corresponding to the best achievable solution in some class of learned functions.

(Bengio et al. Deep Learning 2014)



# Exercise question



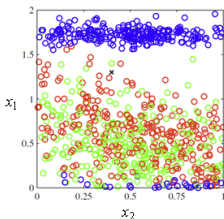
- In the figure above, the training data (10 black dots) were selected from a quadratic function plus Gaussian noise, i.e.,  $f(x) = w_2x^2 + w_1x + b + \epsilon$  where  $p(\epsilon) = N(0, \sigma^2)$ . The degree-10 polynomial fits the data perfectly. Which learner should be chosen in order to better predict new examples? The second-order function or the 10th degree function?
- If the ten training examples were generated from a 10th degree polynomial plus Gaussian noise, which learner should be chosen?
- If the one million training examples were generated from a quadratic function plus Gaussian noise, which learner should be chosen?

# How to reduce capacity?

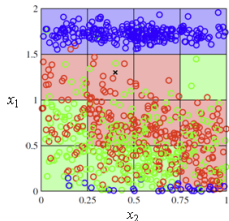
- Reduce the number of features
- Reduce the number of **independent** parameters
- Reduce the network size of deep models
- Reduce the number of training iterations
- Add regularization to the learner
- ...

# Curse of dimensionality

- Why do we need to reduce the dimensionality of the feature space?



Scatter plot of the training data of three classes. Two features are used. The goal is to classify the new testing point denoted by 'x'.

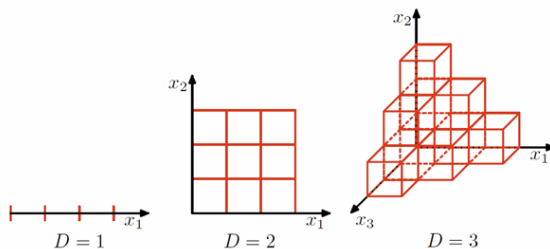


The feature space is uniformly divided into cells. A cell is labeled as a class, if the majority of training examples in that cell are from that class. The testing point is classified according to the label of the cell where it falls in.

(Duda et al. Pattern Classification 2000)

# Curse of dimensionality

- The more training samples in each cell, the more robust the classifier
- The number of cells grows exponentially with the dimensionality of the feature space. If each dimension is divided into three intervals, the number of cells is  $N = 3^D$
- Some cells are empty when the number of cells is very large!



(Duda et al. Pattern Classification 2000)

# Regularization

- Equivalent to imposing a preference over the set of functions that a learner can obtain as a solution
- In Bayesian learning, it is reflected as a prior probability distribution over the space of functions (or equivalently their parameters) that the learner can assess
- Regularization prevents overfitting by adding penalty for complexity
- Training a classifier/regressor is to minimize  
**Prediction error on the training set + regularization**
- Examples

- The objective function for linear regression becomes

$$\text{MSE}_{\text{train}} + \text{regularization} = \frac{1}{N} \sum_i (\mathbf{w}^t \mathbf{x}_i^{(\text{train})} - y_i^{(\text{train})})^2 + \lambda \|\mathbf{w}\|_2^2$$

- Multi-task learning, transfer learning, dropout, sparsity, pre-training

# Function estimation

- We are interested in predicting  $y$  from input  $\mathbf{x}$  and assume there exists a function that describes the relationship between  $y$  and  $\mathbf{x}$ , e.g.  $y = f(\mathbf{x}) + \epsilon$ , where  $\epsilon$  is random noise following certain distribution.
- Prediction function  $f$  can be parametrized by a parameter vector  $\theta$ .
- Estimating  $\hat{f}_n$  from a training set  $\mathcal{D}_n = \{(\mathbf{x}_1^{(\text{train})}, y_1^{(\text{train})}), \dots, (\mathbf{x}_n^{(\text{train})}, y_n^{(\text{train})})\}$  is equivalent to estimating  $\hat{\theta}_n$  from  $\mathcal{D}_n$ .
- Since  $\mathcal{D}_n$  is randomly generated from a underlying distribution, both  $\hat{\theta}$  and  $\hat{f}$  are random variables (or vectors, or functions) distributed according to some probability distributions.
- The quality of estimation can be measured by bias and variance compared with the “true” parameter vector  $\theta$  or function  $\hat{f}$
- With a better design of the parametric form of the function, the learner could achieve low generalization error even with small capacity
- This design process typical involves domain knowledge

$$\text{bias}(\hat{\theta}) = E(\hat{\theta}) - \theta$$

where expectation is over all the train sets of size  $n$  sampled from the underlying distribution

- An estimator is called unbiased if  $E(\hat{\theta}) = \theta$
- Example: Gaussian distribution.  $p(\mathbf{x}_i; \theta) = \mathcal{N}(\theta, \Sigma)$  and the estimator is  $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^{(\text{train})}$

$$E(\hat{\theta}) = E \left[ \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^{(\text{train})} \right] = \frac{1}{n} \sum_{i=1}^n E \left[ \mathbf{x}_i^{(\text{train})} \right] = \frac{1}{n} \sum_{i=1}^n \theta = \theta$$

$$\text{Var}[\hat{\theta}] = E[(\hat{\theta} - E[\hat{\theta}])^2] = E[\hat{\theta}^2] - E[\hat{\theta}]^2$$

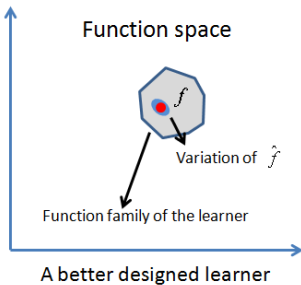
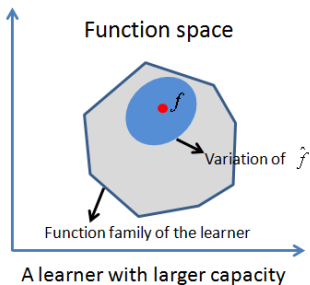
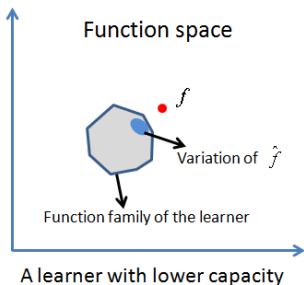
- Variance typically decreases as the size of the train set increases  $\text{Var}(\hat{\theta}_m) = \frac{1}{m}\theta(1-\theta)$  (page 129)

- Both bias and variance are the sources of estimation errors

$$\text{MSE} = E[(\hat{\theta} - \theta)^2] = \text{Bias}(\hat{\theta})^2 + \text{Var}[\hat{\theta}]$$

- Increasing the capacity of a learner may also increase variance, although it has better chance to cover the true function

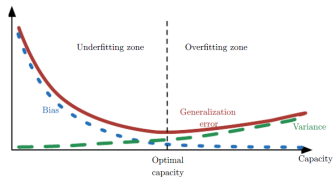




# Summary: issues to be concerned in machine learning

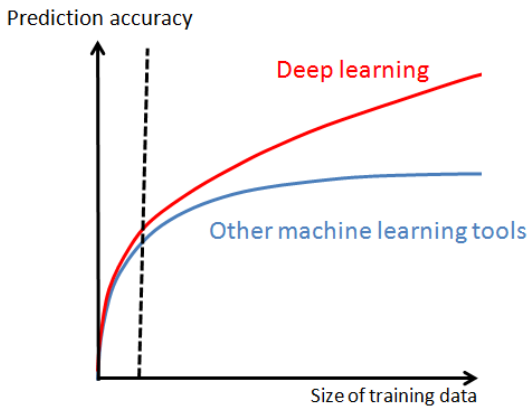
- Effective optimization methods and models to address the underfitting problem
- How to balance the trade-off between capacity and generalization?
- How to effectively reduce capacity (which means also reducing estimation variance) without increasing the bias much?
- For machine learning with big training data, how to effectively increase capacity to cover or get closer to the true function to be estimated?

page 130



# Open discussion

- Why does deep learning have different behavior than other machine learning methods for large scale training?



# Discriminative model

- Directly model  $P(y|\mathbf{x})$  and decision boundaries
- Learn the discriminative functions  $g_k(\mathbf{x})$

$$y = \arg \max_k g_k(\mathbf{x})$$

- In the linear case,  $g_k(\mathbf{x}) = \mathbf{w}_k^t \mathbf{x}$
- $P(y|\mathbf{x})$  can be estimated from the linear discriminant functions

$$P(y = j|\mathbf{x}) = \frac{e^{\mathbf{w}_j^t \mathbf{x}}}{\sum_{k=1}^K e^{\mathbf{w}_k^t \mathbf{x}}}$$

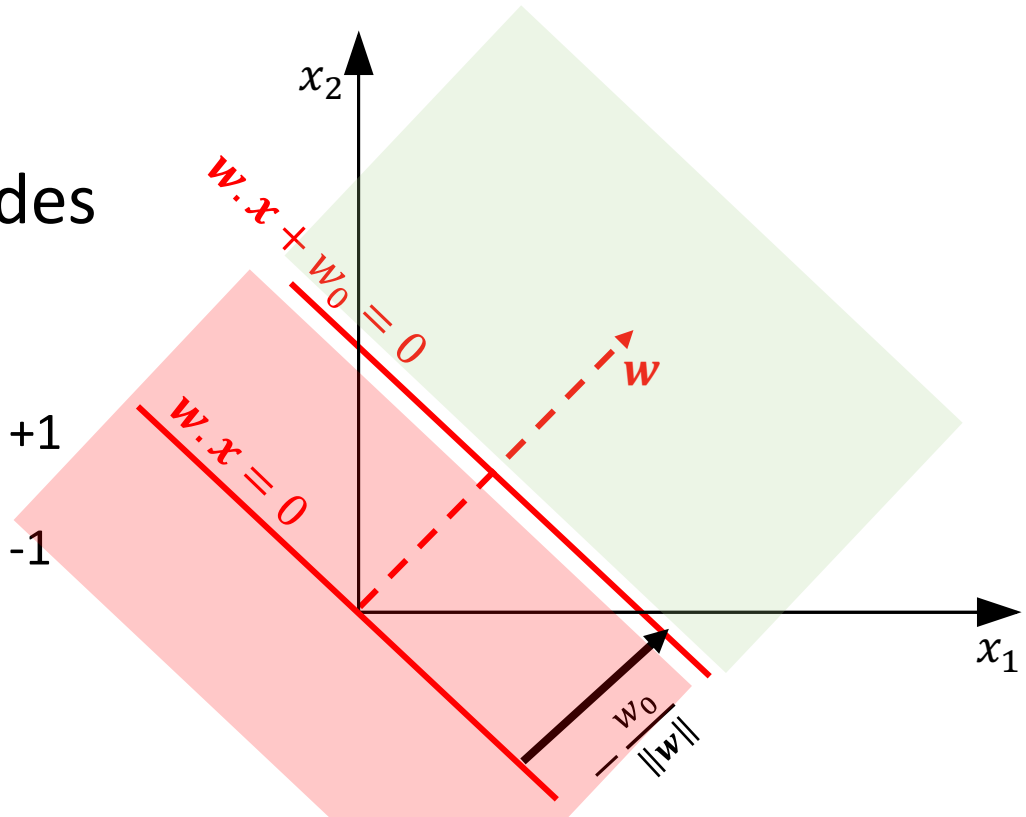
It is also called softmax function

- Examples: SVM, boosting, K-nearest-neighbor

# Linear decision boundaries

$$\hat{y}(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + w_0)$$

- $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 = 0\}$  is a hyperplane in  $\mathbb{R}^d$ 
  - decision boundary
  - $\mathbf{w}$  is direction of normal
  - $w_0$  is the offset
- $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 = 0\}$  divides  $\mathbb{R}^d$  into two halfspaces (regions)
  - $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 \geq 0\}$  get label +1 and  
 $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 < 0\}$  get label -1



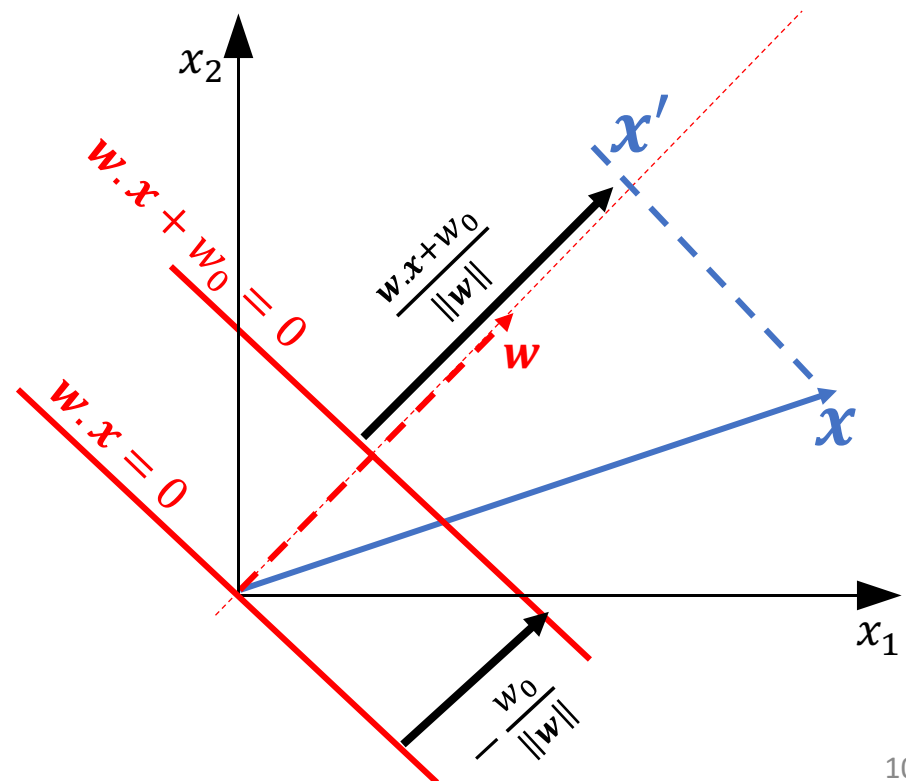
# Linear decision boundaries

$$\hat{y}(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + w_0)$$

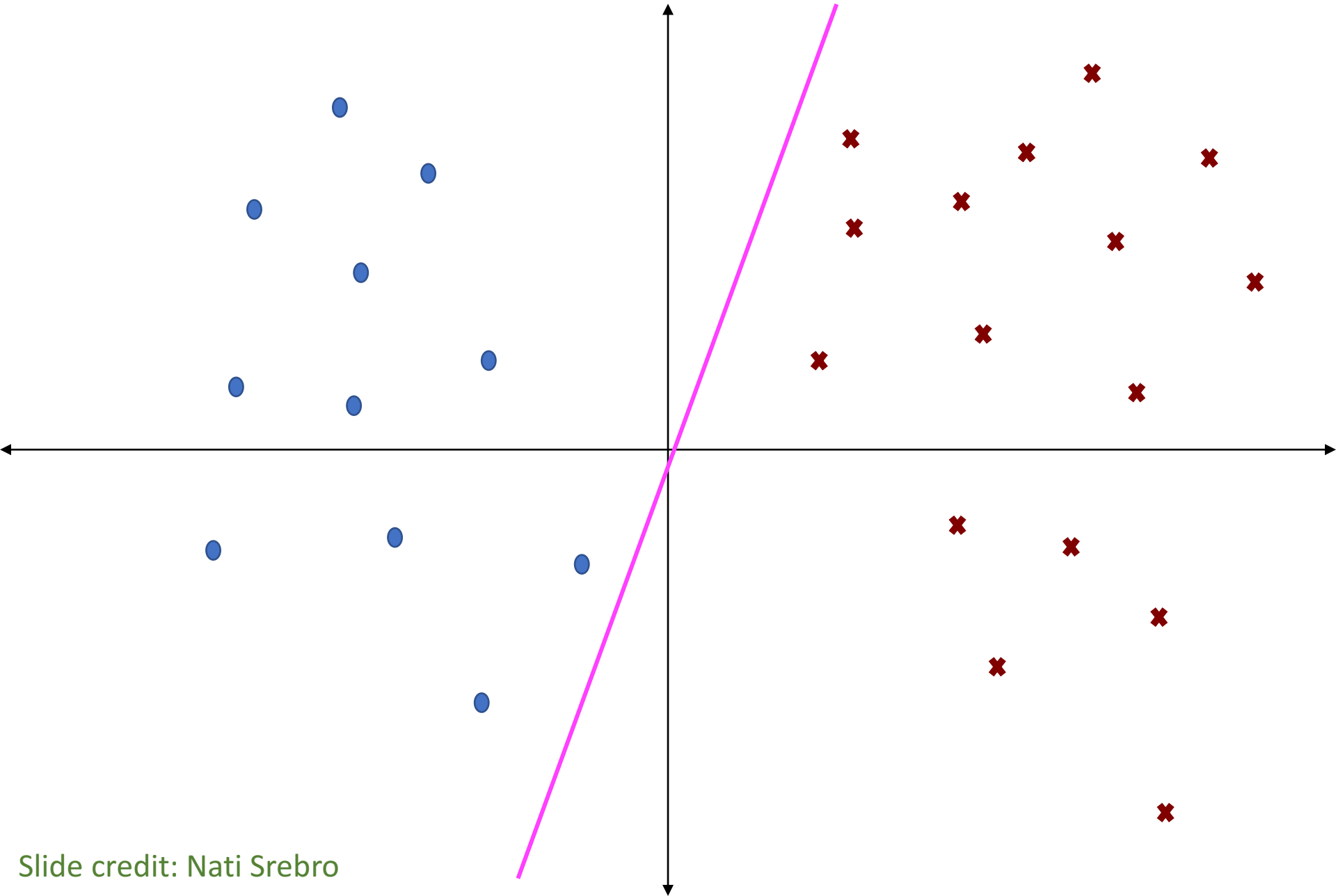
- $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 = 0\}$  is a hyperplane in  $\mathbb{R}^d$ 
  - decision boundary
  - $\mathbf{w}$  is direction of normal
  - $w_0$  is the offset
- $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 = 0\}$  divides  $\mathbb{R}^d$  into two halfspaces (regions)
  - $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 \geq 0\}$  get label +1 and
  - $\{\mathbf{x}: \mathbf{w} \cdot \mathbf{x} + w_0 < 0\}$  get label -1

Maps  $\mathbf{x}$  to a 1D coordinate

$$x' = \frac{\mathbf{w} \cdot \mathbf{x} + w_0}{\|\mathbf{w}\|}$$



# Linear separators in 2D

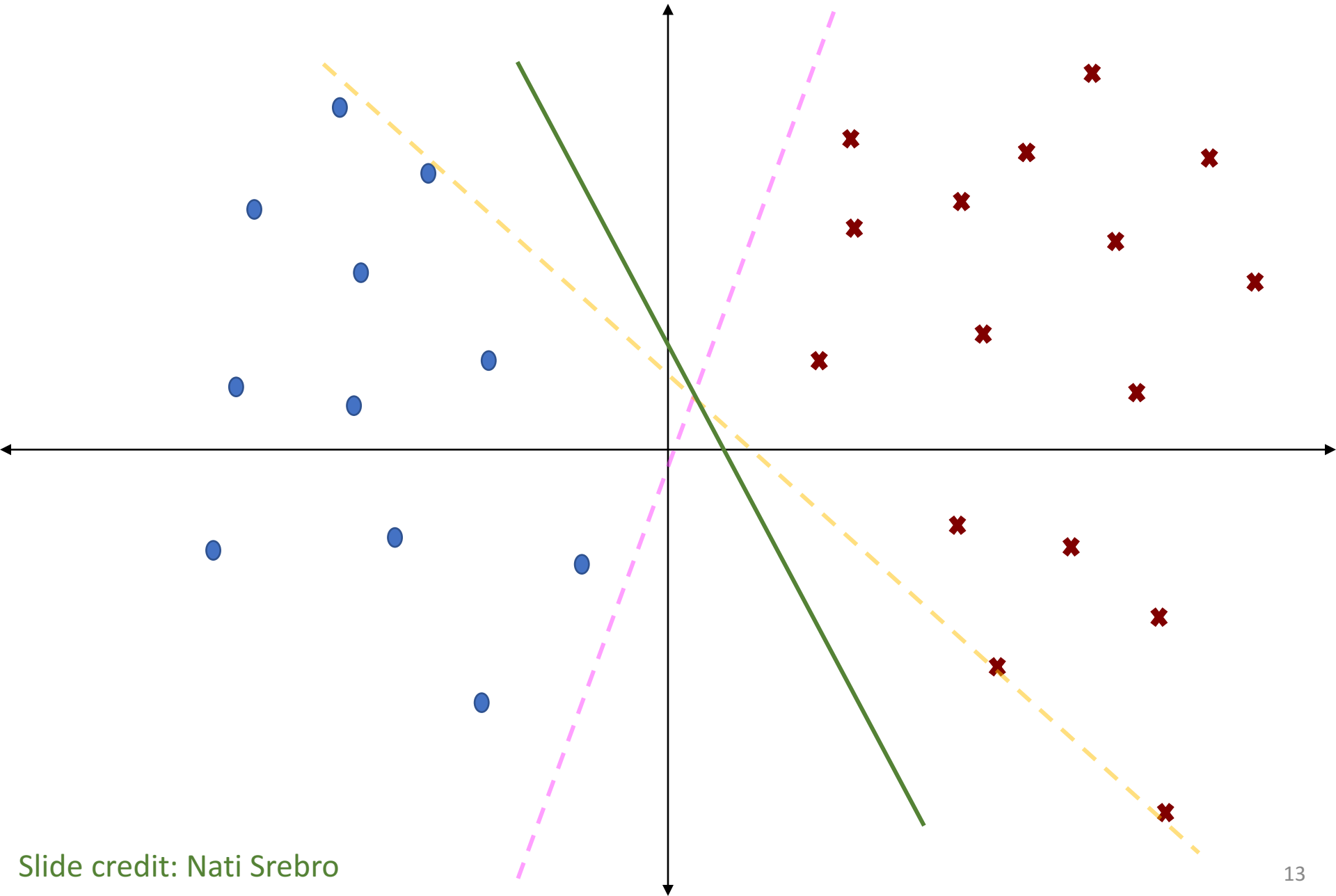


The scatter plot displays two classes of data points: blue circles and red crosses. A solid yellow line represents the linear regression fit, and a dashed pink line represents the decision boundary. The axes are labeled with arrows.

Class	Point ID	X	Y
Blue Circles	1	-1.5	1.5
	2	-1.0	2.5
	3	-0.5	3.5
	4	-0.5	2.0
	5	-0.5	1.0
	6	0.0	3.0
	7	0.5	2.0
	8	1.0	1.0
	9	1.5	0.5
	10	2.0	0.0
	11	2.5	-0.5
	12	3.0	-1.0
Red Crosses	1	0.5	3.0
	2	1.0	3.5
	3	1.5	3.0
	4	2.0	3.5
	5	2.5	3.0
	6	3.0	3.5
	7	3.5	3.0
	8	4.0	3.5
	9	4.5	3.0
	10	5.0	3.5
	11	5.5	3.0
	12	6.0	3.5
	13	6.5	3.0
	14	7.0	3.5
	15	7.5	3.0

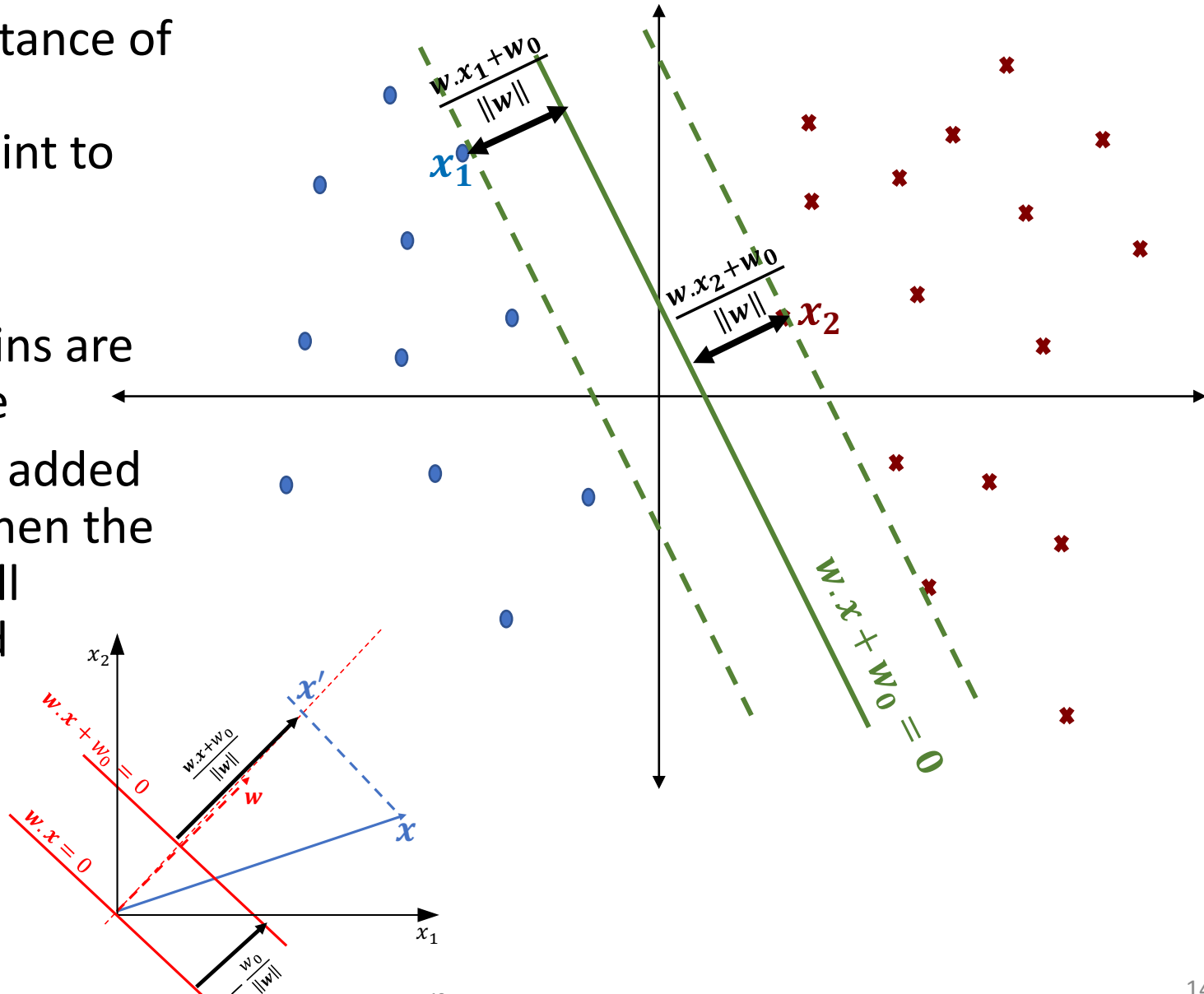


# Linear separators in 2D



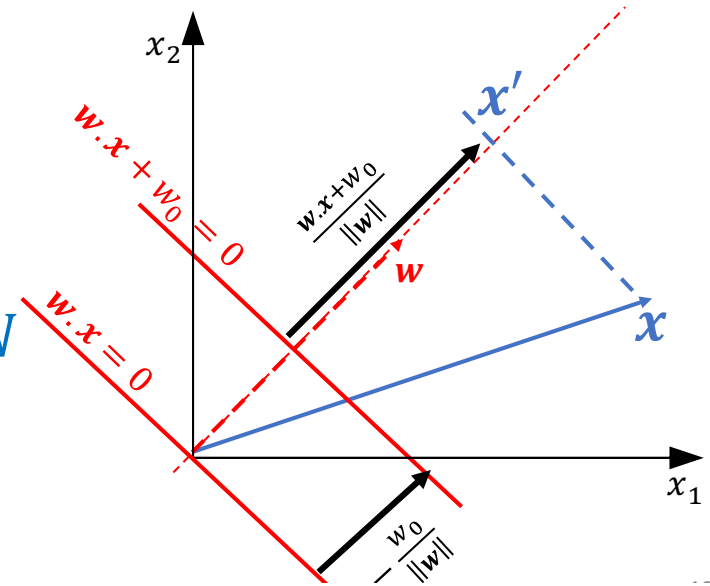
# Margin of a classifier

- **Margin:** distance of the closest instance point to the linear hyperplane
- Large margins are more stable
  - If noise is added to data, then the data is still separated



# Maximum margin classifier

- $S = \{(\mathbf{x}^{(i)}, y^{(i)}): i = 1, 2, \dots, N\}$   
binary classes  $\mathcal{Y} = \{-1, 1\}$
- Assume data is “linearly separable”
  - $\exists \mathbf{w}, w_0$  such that for all  $i = 1, 2, \dots, N$   
 $y^{(i)} = \text{sign}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)$   
 $\Rightarrow y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0) > 0$
- **Maximum margin** separator given by



$$\hat{\mathbf{w}}, \hat{w}_0 = \underset{\mathbf{w} \in \mathbb{R}^d, w_0}{\operatorname{argmax}} \quad \min_i \underbrace{\frac{y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)}{\|\mathbf{w}\|}}_{\text{margin of sample } i}$$

smallest  
margin

margin of  
sample  $i$

# Maximum margin classifier

$$\hat{\mathbf{w}}, \hat{w}_0 = \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^d, w_0 \in \mathbb{R}} \min_i \frac{y^{(i)} (\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)}{\|\mathbf{w}\|}$$

- **Claim 1:** If  $\hat{\mathbf{w}}, \hat{w}_0$  is a solution, then for any  $\gamma > 0$ ,  $\gamma \hat{\mathbf{w}}, \gamma \hat{w}_0$  is also a solution
- **Option 1:** We can fix  $\|\mathbf{w}\| = 1$  to get

$$\hat{\mathbf{w}}, \hat{w}_0 = \operatorname{argmax}_{\|\mathbf{w}\|=1, w_0} \min_i y^{(i)} (\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)$$

# Maximum margin classifier

$$\hat{\mathbf{w}}, \hat{w}_0 = \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^d, w_0 \in \mathbb{R}} \min_i \frac{y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)}{\|\mathbf{w}\|}$$

- **Claim 1:** If  $\hat{\mathbf{w}}, \hat{w}_0$  is a solution, then for any  $\gamma > 0$ ,  $\gamma \hat{\mathbf{w}}, \gamma \hat{w}_0$  is also a solution
- **Option 1:** we can fix  $\|\mathbf{w}\| = 1$  to get

$$\hat{\mathbf{w}}, \hat{w}_0 = \operatorname{argmax}_{\|\mathbf{w}\|=1} \min_i y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)$$

- **Option 2:** we can also fix  $\min_i y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0) = 1$ 
  - margin now is  $\frac{1}{\|\mathbf{w}\|}$
  - Instead of “increasing margin” we can “reduce norm”

# Max-margin classifier equivalent formulation

- Solve:  $\tilde{\mathbf{w}}, \tilde{w}_0 = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w}\|^2$   
s.t.  $\forall i, y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0) \geq 1$

Hard margin  
Support Vector  
Machine (SVM)

- **Claim 2:** Equivalent to previous slide

→  $\frac{\tilde{\mathbf{w}}}{\|\tilde{\mathbf{w}}\|}, \frac{\tilde{w}_0}{\|\tilde{\mathbf{w}}\|}$  is solution for

$$\hat{\mathbf{w}}, \hat{w}_0 = \max_{\|\mathbf{w}\|=1} \min_i y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)$$

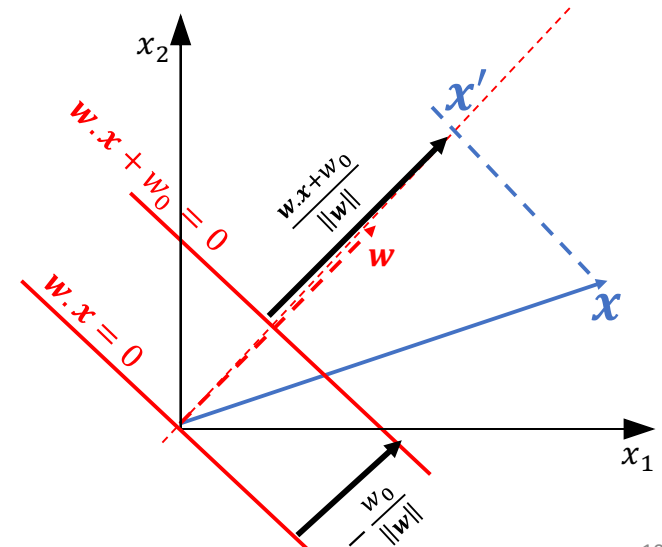
- Proof:

1. Let  $\min_i y^{(i)}(\hat{\mathbf{w}} \cdot \mathbf{x}^{(i)} + \hat{w}_0) = \hat{\gamma}$ , then  $\min_i y^{(i)} \left( \frac{\hat{\mathbf{w}}}{\hat{\gamma}} \cdot \mathbf{x}^{(i)} + \frac{\hat{w}_0}{\hat{\gamma}} \right) \geq 1$
2.  $\Rightarrow \|\tilde{\mathbf{w}}\| \leq \left\| \frac{\hat{\mathbf{w}}}{\hat{\gamma}} \right\| = \frac{1}{\hat{\gamma}}$
3.  $\min_i y^{(i)} \left( \frac{\tilde{\mathbf{w}}}{\|\tilde{\mathbf{w}}\|} \cdot \mathbf{x}^{(i)} + \frac{\tilde{w}_0}{\|\tilde{\mathbf{w}}\|} \right) = \min_i \frac{y^{(i)}(\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)} + \tilde{w}_0)}{\|\tilde{\mathbf{w}}\|} \geq \frac{1}{\|\tilde{\mathbf{w}}\|} \geq \hat{\gamma}$

# Maximum margin classifier formulations

- Original formulation

$$\hat{\mathbf{w}}, \hat{w}_0 = \underset{\mathbf{w} \in \mathbb{R}^d, w_0 \in \mathbb{R}}{\operatorname{argmax}} \min_i \frac{y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0)}{\|\mathbf{w}\|}$$



- Fixing  $\|\mathbf{w}\| = 1$

$$\hat{\mathbf{w}}, \hat{w}_0 = \underset{\mathbf{w}, w_0}{\operatorname{argmax}} \min_i y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0) \quad \text{s.t.} \quad \|\mathbf{w}\| = 1$$

- Fixing  $\min_i y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0) = 1$

$$\tilde{\mathbf{w}}, \tilde{w}_0 = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad \forall i, y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + w_0) \geq 1$$

Henceforth,  $w_0$  will be absorbed into  $\mathbf{w}$  by adding an additional feature of '1' to  $\mathbf{x}$

# Margin and norm

- $\text{margin}(\mathbf{w}) = \min_i \frac{y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)}}{\|\mathbf{w}\|}$
- Remember in regression: small norm solutions have low complexity!
  - Is this true for maximum margin classifiers?
  - what about classification with logistic loss  $\sum_i \log(1 + \exp(-y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)}))$ ?
  - how to do capacity control in maximum margin classifier learning?
- Some places  $\min_i y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)}$  referred as margin  $\rightarrow$  implicitly assumes normalization
  - $\min_i y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)}$  is meaningless without knowing what  $\|\mathbf{w}\|$  is!



# Introduce Slack Variables

## Non-Separable Data

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d} \quad & ||\mathbf{w}||^2 \\ \text{subject to} \quad & y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \geq 1 \quad \forall i = 1 \dots n \end{aligned}$$

$\Downarrow$

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, \xi_1 \dots \xi_n \in \mathbb{R}} \quad & ||\mathbf{w}||^2 + \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \geq 1 - \xi_i \quad \forall i = 1 \dots n \\ & \xi_i \geq 0 \quad \forall i = 1 \dots n \end{aligned}$$

# Unconstrained Formulation

“Soft” SVM solution

$$\mathbf{w}_S^{\text{soft}}, \xi_1^* \dots \xi_n^* := \arg \min_{\mathbf{w} \in \mathbb{R}^d, \xi_1 \dots \xi_n \in \mathbb{R}} \|\mathbf{w}\|^2 + \sum_{i=1}^n \xi_i$$

with constraints  $\xi_i \geq \max(0, 1 - y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)})$  for all  $i = 1 \dots n$ .

Note that  $\xi_i^* = \max(0, 1 - y^{(i)} \mathbf{w}_S^{\text{soft}} \cdot \mathbf{x}^{(i)})$ , so

$$\mathbf{w}_S^{\text{soft}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{w}\|^2 + \sum_{i=1}^n \max(0, 1 - y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)})$$

No constraints :) Convex but not differentiable, can still be optimized by subgradient descent

# Soft SVMs as Empirical Risk Minimization

- ▶ What we really want: minimize the 0-1 loss

$$\mathbf{w}_{\mathcal{S}}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n \left[ \left[ y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \leq 0 \right] \right]$$

where  $[[\tau]]$  is 1 if  $\tau$  is true and 0 otherwise

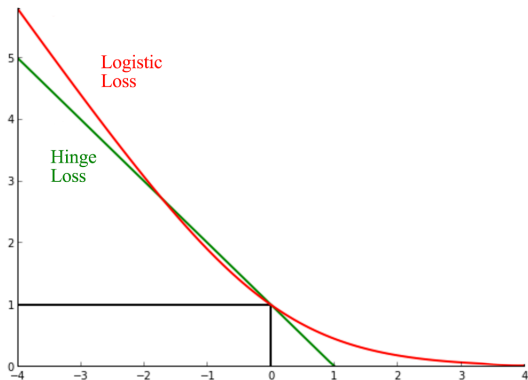
- ▶ Difficult to optimize (neither convex nor differentiable)
- ▶ Instead minimize the **hinge loss**

$$\mathbf{w}_{\mathcal{S}}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n \underbrace{\max \left( 0, 1 - y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \right)}_{\text{hinge}(\mathbf{w} \cdot \mathbf{x}^{(i)})}$$

which is a convex upper bound on the 0-1 loss

# A Big Picture of Binary Classification

Both logistic regression and soft SVMs are  $l_2$ -regularized minimization of the 0-1 loss by convex surrogates.



## Generalized Representer Theorem

**Claim.** Let  $l : \mathbb{R}^n \rightarrow [0, \infty)$  be any function and define

$$\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{w}\|^2 + l\left(\mathbf{w} \cdot \mathbf{x}^{(1)}, \dots, \mathbf{w} \cdot \mathbf{x}^{(n)}\right)$$

Then  $\mathbf{w}^* = \sum_{i=1}^n \beta_i \mathbf{x}^{(i)}$  for some  $\beta_1 \dots \beta_n \in \mathbb{R}$ .

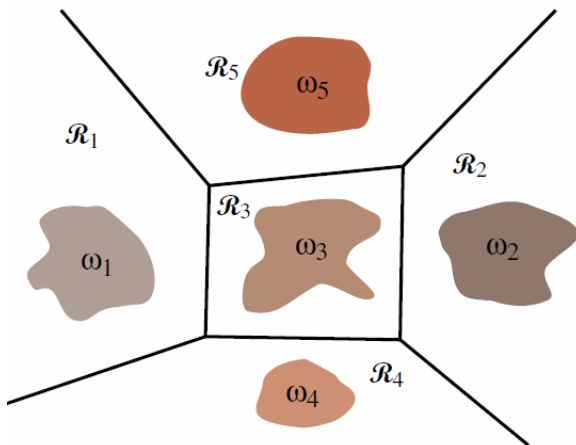
**Proof.** Same as in the hard SVM case

Thus we can similarly derive an inner product formulation of the soft SVM solution (will come back to this later):

$$\mathbf{w}_S^{\text{soft}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{w}\|^2 + \underbrace{\sum_{i=1}^n \max\left(0, 1 - y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)}\right)}_{l(\mathbf{w} \cdot \mathbf{x}^{(1)}, \dots, \mathbf{w} \cdot \mathbf{x}^{(n)})}$$

# Discriminative model

- It is easier for discriminative models to fit data



# Discriminative model

- Parameter  $\theta = \{\mathbf{w}_k\}$  can be estimated from maximizing the data likelihood

$$\hat{\theta} = \arg \max_{\theta} P(\mathcal{D}_n | \theta) = \arg \max_{\theta} \prod_{i=1}^n P(y_n^{(\text{train})} | \mathbf{x}_n^{(\text{train})}, \theta)$$

- Maximum a posteriori* (MAP) estimation

$$\theta = \arg \max_{\theta} p(\theta | \mathcal{D}_n) = \arg \max_{\theta} \log P(\mathcal{D}_n | \theta) + \log p(\theta)$$

- According to the Bayes' rule, i.e.,  $p(\theta | \mathcal{D}_n) = P(\mathcal{D}_n | \theta)p(\theta)/P(\mathcal{D}_n)$ ,

$$\theta = \arg \max_{\theta} \log P(\mathcal{D}_n | \theta) + \log p(\theta)$$

$$\theta = \arg \max_{\theta} \sum_{i=1}^n \log P(y_n^{(\text{train})} | \mathbf{x}_n^{(\text{train})}, \theta) + \log p(\theta)$$

- prior  $p(\theta)$  corresponds to a regularizer, e.g.

$$p(\theta) = e^{-\lambda \|\theta\|^2}$$

# Generative model

- Estimate the underlying class conditional probability densities  $p(\mathbf{x}|y = k)$  and then construct the classifier using the Bayesian decision theory

$$P(y = k|\mathbf{x}) = \frac{p(\mathbf{x}|y = k)P(y = k)}{p(\mathbf{x})} = \frac{p(\mathbf{x}|y = k)P(y = k)}{\sum_{k'=1}^K p(\mathbf{x}|y = k')P(y = k')}$$

- $p(\mathbf{x}|y)$  and  $P(y)$  are parameterized by  $\theta$
- Prior  $P(y)$  can be used to model the dependency among predictions, such as the segmentation labels of pixels or predictions of speech sequences.
- It is more difficult to model class conditional probability densities. However, it also adds stronger regularization to model fitting, since the learned model not only needs to predict class labels but also generate the input data.
- It is easier to add domain knowledge when designing the models of  $p(\mathbf{x}|y)$



# Supervised and unsupervised learning

- Supervised learning: the goal is to use input-label pairs,  $(\mathbf{x}, y)$  to learn a function  $f$  that predicts a label (or a distribution over labels) given the input,  $\hat{y} = f(\mathbf{x})$
- Unsupervised learning: no label or other target is provided. The data consists of a set of examples  $\mathbf{x}$  and the objective is to learn about the statistical structure of  $\mathbf{x}$  itself.
- Weakly supervised learning: the training data contains  $(\mathbf{x}, y)$  pairs as in supervised learning, but the labels  $y$  are either unreliably present (i.e. with missing values) or noisy (i.e. where the label given is not the true label)

# Unsupervised learning

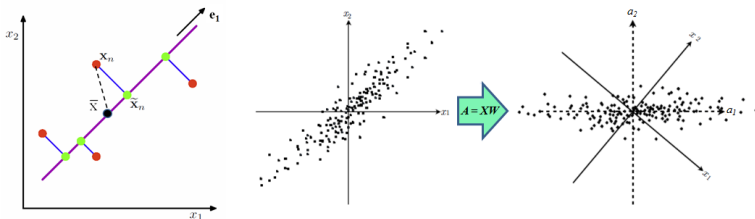
- Find the “best” representation of data that reserves as much information about  $\mathbf{x}$  as possible while being “simpler” than  $\mathbf{x}$   
Taking linear case as an example

$$\tilde{\mathbf{x}} = a_0 + \sum_{i=1}^{d'} a_i \mathbf{e}_i$$

- Lower dimensional representation:  $d' < d$
  - Sparse representation: the number of non-zero  $a_i$  is small
  - Independent representation: disentangle the sources of variations underlying the data distributions such that the dimensions of the representation are statistically independent, i.e.  $a_i$  and  $a_j$  are statistically independent
- Deep learning is to learn data representation, but in a nonlinear and hierarchical way

# Principal Component Analysis (PCA)

- There are  $n$   $d$ -dimensional samples  $\mathbf{x}_1, \dots, \mathbf{x}_n$ .
- PCA seeks a **principal subspace** spanned by  $d'$  ( $d' < d$ ) orthonormal vectors  $\mathbf{e}_1, \dots, \mathbf{e}_{d'}$ , such that
  - the projected samples ( $\tilde{\mathbf{x}}_k = a_0 + \sum_{i=1}^{d'} a_{ki} \mathbf{e}_i$ ) onto this subspace has maximum variance; or equivalently
  - the mean squared distance between the samples and their projections are minimized.
  - The projections  $\{a_{ki}\}$  are uncorrelated (i.e. independent if data distribution is assumed as Gaussian distribution)



# Formulation of PCA

$$\arg \max_{\{\mathbf{e}_i\}} \sum_{k=1}^n \|\tilde{\mathbf{x}}_k - \bar{\mathbf{x}}\|^2$$

or

$$\arg \min_{\{\mathbf{e}_i\}} \sum_{k=1}^n \|\mathbf{x}_k - \tilde{\mathbf{x}}_k\|^2$$

$$(\|\mathbf{x}_k - \bar{\mathbf{x}}\|^2 = \|\mathbf{x}_k - \tilde{\mathbf{x}}_k\|^2 + \|\tilde{\mathbf{x}}_k - \bar{\mathbf{x}}\|^2)$$

# Zero- and One-Dimensional Representations by PCA

- zero-dimensional representation: use a single vector  $\mathbf{x}_0$  to represent all the samples and minimize the squared-error function

$$J_0(\mathbf{x}_0) = \sum_{k=1}^n \|\mathbf{x}_0 - \mathbf{x}_k\|^2$$

The solution is  $\mathbf{x}_0 = \bar{\mathbf{x}} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$

- One-dimensional representation: project the data to a line running through the sample mean,  $\tilde{\mathbf{x}}_k = \bar{\mathbf{x}} + a_{k1}\mathbf{e}_1$  and minimize the squared-error criterion function

$$J_1(a_{11}, \dots, a_{n1}, \mathbf{e}_1) = \sum_{k=1}^n \|\tilde{\mathbf{x}}_k - \mathbf{x}_k\|^2$$

# Find the Principal Components $a_{k1}$

$$\begin{aligned} & J_1(a_{11}, \dots, a_{n1}, \mathbf{e}_1) \\ &= \sum_{k=1}^n \|(\bar{\mathbf{x}} + a_{k1}\mathbf{e}_1) - \mathbf{x}_k\|^2 = \sum_{k=1}^n \|a_{k1}\mathbf{e}_1 - (\mathbf{x}_k - \bar{\mathbf{x}})\|^2 \\ &= \sum_{k=1}^n a_{k1}^2 \|\mathbf{e}_1\|^2 - 2 \sum_{k=1}^n a_{k1} \mathbf{e}_1^t (\mathbf{x}_k - \bar{\mathbf{x}}) + \sum_{k=1}^n \|\mathbf{x}_k - \bar{\mathbf{x}}\|^2 \end{aligned}$$

Since  $\mathbf{e}_1$  is a unit vector,  $\|\mathbf{e}_1\| = 1$ . To minimize  $J_1$ , set  $\frac{\partial J_1}{\partial a_{k1}} = 0$  and we have

$$a_{k1} = \mathbf{e}_1^t (\mathbf{x}_k - \bar{\mathbf{x}})$$

- We obtain a least-squares solution by projecting the vector  $\mathbf{x}_k$  onto the line in the direction of  $\mathbf{e}_1$  passing through the mean.

# Find the Optimal Projection Direction $\mathbf{e}_1$

$$\begin{aligned} J_1(\mathbf{e}_1) &= \sum_{k=1}^n a_{k1}^2 - 2 \sum_{k=1}^n a_{k1}^2 + \sum_{k=1}^n \|\mathbf{x}_k - \bar{\mathbf{x}}\|^2 \\ &= - \sum_{k=1}^n [\mathbf{e}_1^t (\mathbf{x}_k - \bar{\mathbf{x}})]^2 + \sum_{k=1}^n \|\mathbf{x}_k - \bar{\mathbf{x}}\|^2 \\ &= - \sum_{k=1}^n \mathbf{e}_1^t (\mathbf{x}_k - \bar{\mathbf{x}}) (\mathbf{x}_k - \bar{\mathbf{x}})^t \mathbf{e}_1 + \sum_{k=1}^n \|\mathbf{x}_k - \bar{\mathbf{x}}\|^2 \\ &= - \mathbf{e}_1^t \mathbf{S} \mathbf{e}_1 + \sum_{k=1}^n \|\mathbf{x}_k - \bar{\mathbf{x}}\|^2 \end{aligned}$$

# Find the Optimal Projection Direction $\mathbf{e}_1$

- $\mathbf{S} = \sum_{k=1}^n (\mathbf{x}_k - \bar{\mathbf{x}})(\mathbf{x}_k - \bar{\mathbf{x}})^t$  is the scatter matrix
- $\mathbf{e}_1^t \mathbf{S} \mathbf{e}_1 = \sum_{k=1}^n a_{k1}^2$  ( $\sum_{k=1}^n a_{k1} = 0$ ) is the variance of the projected data
- The vector  $\mathbf{e}_1$  that minimizes  $J_1$  also maximizes  $\mathbf{e}_1^t \mathbf{S} \mathbf{e}_1$ , subject to the constraint that  $\|\mathbf{e}_1\| = 1$



# Lagrange Optimization

- Seek the position  $\mathbf{x}_0$  of an extremum of a scalar-valued function  $f(\mathbf{x})$  subject to the constrain that  $g(\mathbf{x}) = 0$
- First from the Lagrangian function

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

where  $\lambda$  is a scalar called the Lagrange *undetermined multiplier*.

- Convert into an unconstrained problem by taking the derivative,

$$\frac{\partial L(\mathbf{x}, \lambda)}{\partial \mathbf{x}} = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} + \lambda \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} = 0$$

- Solve  $\mathbf{x}$  and  $\lambda$  considering  $g(\mathbf{x}) = 0$

# Find the Optimal Projection Direction $\mathbf{e}_1$

- Use the method of Lagrange multipliers to maximize the  $\mathbf{e}_1^t \mathbf{S} \mathbf{e}_1$  subject to the constraint that  $\|\mathbf{e}_1\| = 1$ ,

$$u = \mathbf{e}_1^t \mathbf{S} \mathbf{e}_1 - \lambda(\mathbf{e}_1^t \mathbf{e}_1 - 1),$$

$$\frac{\partial u}{\partial \mathbf{e}_1} = 2\mathbf{S} \mathbf{e}_1 - 2\lambda \mathbf{e}_1 = 0.$$

- Setting the gradient vector equal to zero, we see that  $\mathbf{e}_1$  must be an eigenvector of the scatter matrix

$$\mathbf{S} \mathbf{e}_1 = \lambda \mathbf{e}_1$$

- Since  $\mathbf{e}_1^t \mathbf{S} \mathbf{e}_1 = \lambda \mathbf{e}_1^t \mathbf{e}_1 = \lambda$ , to maximize  $\mathbf{e}_1^t \mathbf{S} \mathbf{e}_1$ , we select the eigenvector with the largest eigenvalue

# $d'$ -Dimensional Representation by PCA

- $d'$ -dimensional representation:  $\tilde{\mathbf{x}}_k = \bar{\mathbf{x}} + \sum_{i=1}^{d'} a_{ki} \mathbf{e}_i$
- Mean-squared criterion function:

$$J_{d'} = \sum_{k=1}^n \left\| \left( \bar{\mathbf{x}} + \sum_{i=1}^{d'} a_{ki} \mathbf{e}_i \right) - \mathbf{x}_k \right\|^2$$

- Define additional principal components in an incremental fashion by choosing each new direction minimizing  $J$  amongst all possible directions orthogonal to those already considered
- To minimize  $J_{d'}$ ,  $\mathbf{e}_1, \dots, \mathbf{e}_{d'}$  are the  $d'$  eigenvectors of the scatter matrix with the largest eigenvalues.  $a_{ki}$  are the **principal components** of samples.

# $d'$ -Dimensional Representation by PCA

- Since the scatter matrix is real and symmetric, its eigenvectors are orthogonal and its eigenvalues are nonnegative.
- The squared error:

$$J_{d'} = - \sum_{i=1}^{d'} \lambda_i + \sum_{k=1}^n \|\mathbf{x}_k - \bar{\mathbf{x}}\|^2$$

$$J_d = 0 \Rightarrow \sum_{k=1}^n \|\mathbf{x}_k - \bar{\mathbf{x}}\|^2 = \sum_{i=1}^d \lambda_i$$

$$J_{d'} = \sum_{i=d'+1}^d \lambda_i$$

# Variance of Data Captured by the PCA Subspace

- The variance of data projected onto the first  $d'$  eigenvectors is  $\sum_{i=1}^{d'} \lambda_i$
- Measure how much variance has been captured by the first  $d'$  eigenvectors:

$$\frac{\sum_{i=1}^{d'} \lambda_i}{\sum_{j=1}^d \lambda_j}$$

# Covariance Matrix of Principal Components

- The correlation between projections on  $\mathbf{e}_i$  and  $\mathbf{e}_j$  ( $i \neq j$ ) is

$$\sum_{k=1}^n a_{ki} a_{kj} = \sum_{k=1}^n \mathbf{e}_i^t (\mathbf{x}_k - \bar{\mathbf{x}}) (\mathbf{x}_k - \bar{\mathbf{x}})^t \mathbf{e}_j = \mathbf{e}_i^t \mathbf{S} \mathbf{e}_j = \lambda_i \mathbf{e}_i^t \mathbf{e}_j = 0$$

- The covariance matrix of samples in the PCA subspace is  $\text{diag}[\lambda_1, \dots, \lambda_{d'}]$

# Summary of PCA

- The principal subspace is spanned by  $d'$  orthonormal vectors  $\mathbf{e}_1, \dots, \mathbf{e}_{d'}$  which are computed as the  $d'$  eigenvectors of the scatter matrix  $\mathbf{S}$  with the largest eigenvalues  $\lambda_1, \dots, \lambda_{d'}$ .
- The principal components  $a_{ki}$  of the samples are computed as  $a_{ki} = \mathbf{e}_i^t(\mathbf{x}_k - \bar{\mathbf{x}})$
- The variance of the projected samples onto this principal subspace is  $\sum_{i=1}^{d'} \lambda_i$ .
- The mean squared distance between the samples and their projections are  $\sum_{i=d'+1}^d \lambda_i$ .
- PCA disentangles the factors of variation underlying the data, assuming such variation is a Gaussian distribution
- We are interested in learning representations that disentangle more complicated forms of feature dependencies

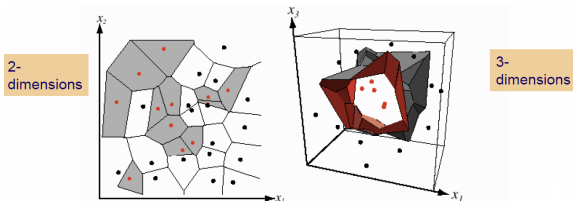
# Smoothness Prior

- Shallow models assume smoothness prior on the prediction function to be learned, i.e.

$$f^*(\mathbf{x}) \approx f^*(\mathbf{x} + \epsilon)$$

where  $\epsilon$  is a small change.

- K-nearest neighbor predictors assume piecewise constant
  - For classification and  $K = 1$ ,  $f(\mathbf{x})$  is the output class associated with the nearest neighbor of  $\mathbf{x}$  in the training set
  - For regression,  $f(\mathbf{x})$  is the average of the outputs associated with the  $K$  nearest neighbors of  $\mathbf{x}$
  - The number of distinguishable regions cannot be more than the number of training examples





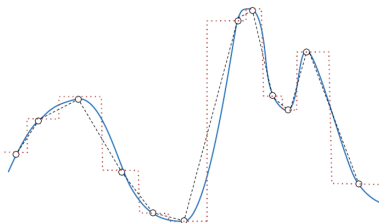
# Interpolation with Kernel

$$f(\mathbf{x}) = b + \sum_{i=1}^n \alpha_i K(\mathbf{x}, \mathbf{x}_i)$$

$K$  is a kernel function, e.g., the Gaussian kernel

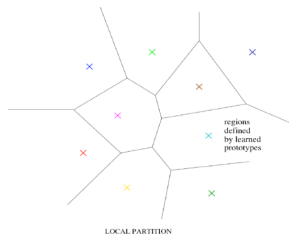
$$K(\mathbf{u}, \mathbf{v}) = N(\mathbf{u} - \mathbf{v}; 0, \sigma^2 I)$$

- $b$  and  $\alpha_i$  can be learned by SVM
- Treat each  $\mathbf{x}_i$  as a template and the kernel function as a similarity function that matches a template and a test example



# Local Representation

- One can think of the training samples as control knots which locally specify the shape of the prediction function
- The smoothness prior only allows the learner to generalize locally. If  $(\mathbf{x}_i, y_i)$  is a supervised training example and  $\mathbf{x}_j$  is a near neighbor of  $\mathbf{x}$ , we expect that  $f(\mathbf{x}) \approx y_j$ . Better generalization can be obtained with more neighbors.
- To distinguish  $O(N)$  regions in the input space, shallow models require  $O(N)$  examples (and typically there are  $O(N)$  parameters associated with the  $O(N)$  regions).



Yoshua Bengio, Ian Goodfellow and Aaron Courville, Chapter 1  
“Machine Learning Basics” in “Deep Learning,” Book in  
preparation for MIT Press, 2014.

<http://www.iro.umontreal.ca/bengioy/dlbook>