

# CSCI 567: Machine Learning

Vatsal Sharan  
Fall 2022

Lecture 3, Sep 8

## Administrivia

HW 1

- HW 2 due in about 1 week (9/14 at 2pm).
- Each person gets 1 late day, if you want to use a late day we'll ask to fill a form (late day counts towards the group member filling form).
- Max 1 late day per HW.
- Please submit your groups by end of today (form on Ed Discussion post by Rachitha).

# Recap

# Supervised learning in one slide

- Loss function:** What is the right loss function for the task?
- Representation:** What class of functions should we use?
- Optimization:** How can we efficiently solve the empirical risk minimization problem?
- Generalization:** Will the predictions of our model transfer gracefully to unseen examples?

*All related! And the fuel which powers everything is **data**.*

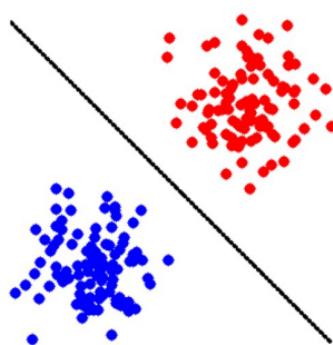
## Summary: Optimization methods

- **GD/SGD** is a first-order optimization method.
- GD/SGM converges to a stationary point. For convex objectives, this is all we need. For nonconvex objectives, it is possible to get stuck at local minimizers or “bad” saddle points (random initialization escapes “good” saddle points).
- **Newton’s method** is a second-order optimization method.
- Newton’s method has a much faster convergence rate, but each iteration also takes much longer. Usually for large scale problems, GD/SGD and their variants are the methods of choice.

# Linear classifiers

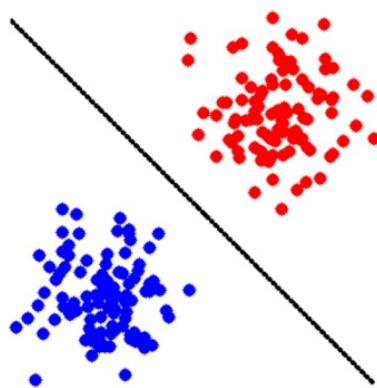
Binary classification:

- input (feature vector):  $x \in \mathbb{R}^d$
- output (label):  $y \in \{-1, +1\}$ .
- goal: learn a mapping  $f : \mathbb{R}^d \rightarrow \{-1, +1\}$



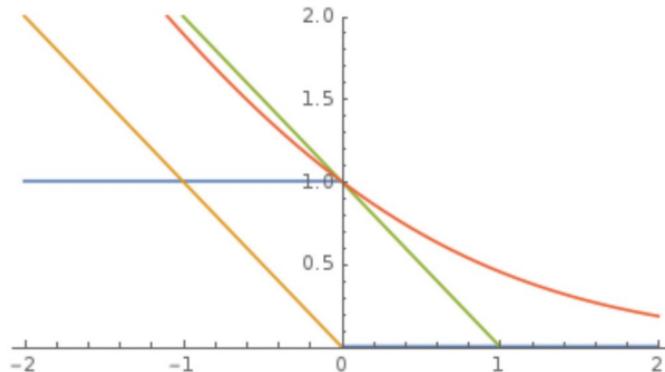
## Representation

Definition: The **function class of separating hyperplanes** is defined as  
 $\mathcal{F} = \{f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x}): \mathbf{w} \in \mathbb{R}^d\}$ .



## Loss function

Use a **convex surrogate loss**



- **perceptron loss**  $\ell_{\text{perceptron}}(z) = \max\{0, -z\}$  (used in Perceptron)
- **hinge loss**  $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\}$  (used in SVM and many others)
- **logistic loss**  $\ell_{\text{logistic}}(z) = \log(1 + \exp(-z))$  (used in logistic regression; the base of log doesn't matter)

## Optimization

Empirical risk minimization (ERM) problem:

$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(y_i \mathbf{w}^T \mathbf{x}_i)$$

Solve using a suitable optimization algorithm:

- **GD:**  $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla F(\mathbf{w})$
- **SGD:**  $\mathbf{w} \leftarrow \mathbf{w} - \eta \tilde{\nabla} F(\mathbf{w})$   $(\mathbb{E}[\tilde{\nabla} F(\mathbf{w})] = \nabla F(\mathbf{w}))$
- **Newton:**  $\mathbf{w} \leftarrow \mathbf{w} - (\nabla^2 F(\mathbf{w}))^{-1} \nabla F(\mathbf{w})$

# Maximum likelihood estimation

*What we observe are labels, not probabilities.*

Take a **probabilistic view**

- assume data is independently generated in this way by some  $w$
- perform Maximum Likelihood Estimation (MLE)

Specifically, what is the probability of seeing labels  $y_1, \dots, y_i$  given  $x_1, \dots, x_i$ , as a function of some  $w$ ?

$$P(w) = \prod_{i=1}^n \mathbb{P}(y_i | x_i; w)$$

**MLE**: find  $w^*$  that **maximizes the probability  $P(w)$**

Minimizing logistic loss is exactly doing MLE for the sigmoid model!



# Generalization

# Reviewing definitions

- Input space:  $\mathcal{X}$
- Output space:  $\mathcal{Y}$
- Predictor:  $f(\mathbf{x}) : \mathcal{X} \rightarrow \mathcal{Y}$
- Distribution  $D$  over  $(\mathbf{x}, y)$ .
- Let  $D^n$  denote the distribution of  $n$  samples  $\{(\mathbf{x}_i, y_i), i \in [n]\}$  drawn i.i.d. from  $D$ .  
*New notation.*
- Risk of a predictor  $f(\mathbf{x})$  is  $R(f) = \mathbb{E}_{(\mathbf{x}, y) \sim D} [\ell(f(\mathbf{x}), y)]$
- Consider the 0-1 loss,  $\ell(f(\mathbf{x}, y)) = \mathbb{1}(f(\mathbf{x}) \neq y)$ .

*The analysis we'll do could also help you solve Problem 3 on HW1.*

## Assumptions for today's theory

Finite sized function classes.

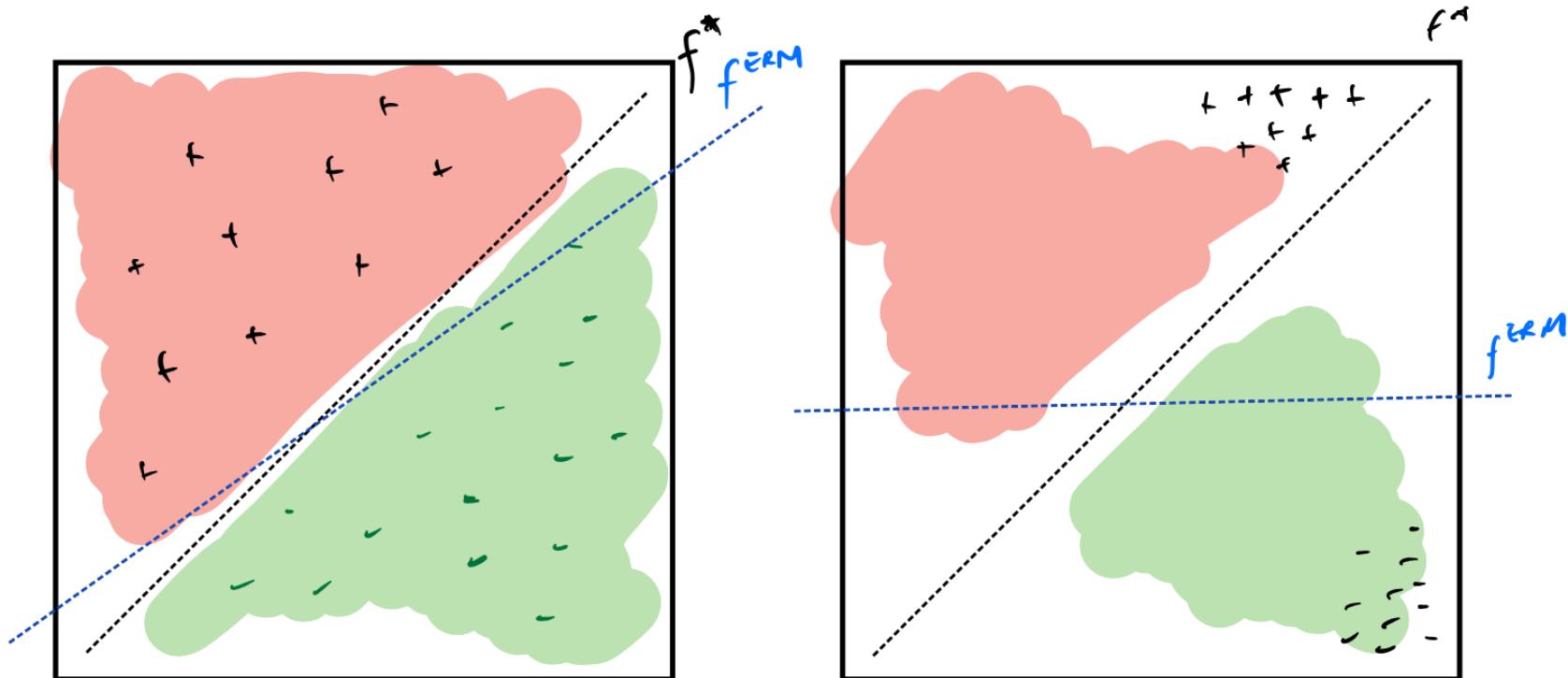
Def: A function class  $\mathcal{F}$  is finite-sized if  $|\mathcal{F}|$  is finite.

e.g.  $\hat{\mathcal{F}} = \{ f(x) = w^T x : w \in \{-1, 0, 1\}^d \}$

$$|\hat{\mathcal{F}}| = 3^d$$

Realizability: There exists  $f^* \in \mathcal{F}$  s.t.  $y = f^*(x) \nabla x \in \mathcal{X}$ .

# Intuition: When does ERM generalize?



Distribution over  $\mathbf{x} = (x_1, x_2)$   
is uniform



Same here

Theorem: Let  $\mathcal{F}$  be a function class with size  $|\mathcal{F}|$ . Let  $y = f^*(x)$  for some  $f \in \mathcal{F}$ . Suppose we get a training set  $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$  of size  $n$  drawn i.i.d. from the data distribution  $D$ . Let

$$f_S^{\text{ERM}} = \underset{f \in \mathcal{F}}{\operatorname{arg\,min}} \frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i).$$

If  $n \geq \frac{\ln(|\mathcal{F}|/\delta)}{\varepsilon^2}$ , then with probability  $(1-\delta)$  over  $\{(x_i, y_i), i \in [n]\}$ ,  $R(f_S^{\text{ERM}}) \leq \varepsilon$  (for constants  $\varepsilon, \delta$ ).

E.g. if  $\Sigma = 0.1$ ,  $\delta = 0.1$ , then with  $n \geq 10 \cdot \ln(10) / \Sigma^2$   
samples, with probability 0.9,  $R(\hat{f}_{\text{ERM}}) \leq 0.1$ .

Proof: Note that there exists  $f^* \in \mathcal{F}$  s.t.  $R(f^*) = 0$ .

$$\mathcal{F}_{\text{bad}} = \{ f \in \mathcal{F} : R(f) > \varepsilon \}$$

Goal ①: What is the probability of "getting tricked"  
by one fixed  $f \in \mathcal{F}_{\text{bad}}$ ?

Consider some  $f' \in \hat{F}_{bad}$ .

$$\Pr_{S \sim D^n} [f' \text{ is an ERM}] \quad (f^* \text{ gets } \hat{R}_S(f^*) = 0)$$

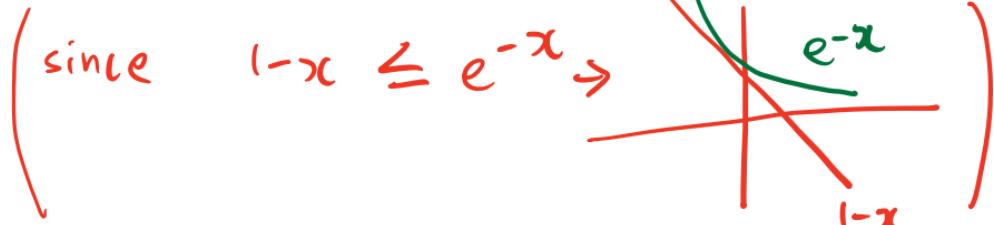
$$= \Pr_{S \sim D^n} \left[ \left\{ \forall i \in \{1, \dots, n\}, f'(x_i) = f^*(x_i) \right\} \right]$$

$$= \prod_{i=1}^n \Pr_{x_i \sim D} [f'(x_i) = f^*(x_i)] \quad \begin{array}{l} \text{(i.i.d. assumption, } \Pr[E_1 \cap E_2] = \Pr(E_1)\Pr(E_2) \text{)} \\ \text{(if } E_1, E_2 \text{ are independent)} \end{array}$$

$$\leq \prod_{i=1}^n (1 - \varepsilon) \quad (f' \in \hat{F}_{bad}, x_i \sim D)$$

$$\leq \prod_{i=1}^n e^{-\varepsilon}$$

$$= e^{-\varepsilon n}.$$



Goal ②: What is the probability of being tricked by any  $f \in \mathcal{F}_{\text{bad}}$ .

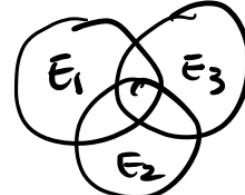
Union bound:

$$\Pr_{S \sim D} \left[ \bigcup_{f \in \mathcal{F}_{\text{bad}}} \{ f \text{ is an ERM} \} \right]$$

$$\leq \sum_{f \in \mathcal{F}_{\text{bad}}} \Pr [ f \text{ is an ERM} ]$$

$$\leq \sum_{f \in \mathcal{F}_{\text{bad}}} e^{-\varepsilon_n} = |\mathcal{F}_{\text{bad}}| e^{-\varepsilon_n} \leq 1 + e^{-\varepsilon_n}.$$

$$\Pr (E_1 \cup E_2 \cup E_3) \leq \sum_{i=1}^3 \Pr (E_i)$$



$$\therefore \text{if } n \geq \frac{1}{\varepsilon} \left( \ln(|\mathcal{F}|) + \ln(1/\delta) \right)$$

then  $\Pr_{S \sim D^n} \left[ \bigcup_{f \in \mathcal{F}_{\text{bad}}} \{ f \text{ is an ERM} \} \right] \leq \delta.$

$$\therefore \text{if } n \geq \frac{\ln(|\mathcal{F}|/\delta)}{\varepsilon}, \text{ w.p. } 1-\delta$$

$$R(f_s^{\text{ERM}}) \leq \varepsilon.$$

Note that in this case the empirical risk  $\hat{R}_s(f_s^{\text{ERM}}) = 0$ , since  $\hat{R}_s(f^*) = 0$  ( $f^* \in \mathcal{F}$  always fits training set perfectly).  $\therefore$  Generalization gap:  $R(f_s^{\text{ERM}})$

## Relaxing our assumptions

- We assumed that the function class is finite-sized. Results can be extended to **infinite function classes** (such as separating hyperplanes).
- We considered 0-1 loss. Can extend to **real-valued loss** (such as for regression).
- We assumed realizability. Can prove similar theorem which guarantees small generalization gap **without realizability** (but with an  $\epsilon^2$  instead of  $\epsilon$  in the denominator). This is called agnostic learning.

## Rule of thumb for generalization

Suppose the functions  $f$  in our function class  $\mathcal{F}$  have  $d$  parameters which can be set. Assume we discretize these parameters so they can take  $W$  possible values each. How much data do we need to have small generalization gap?

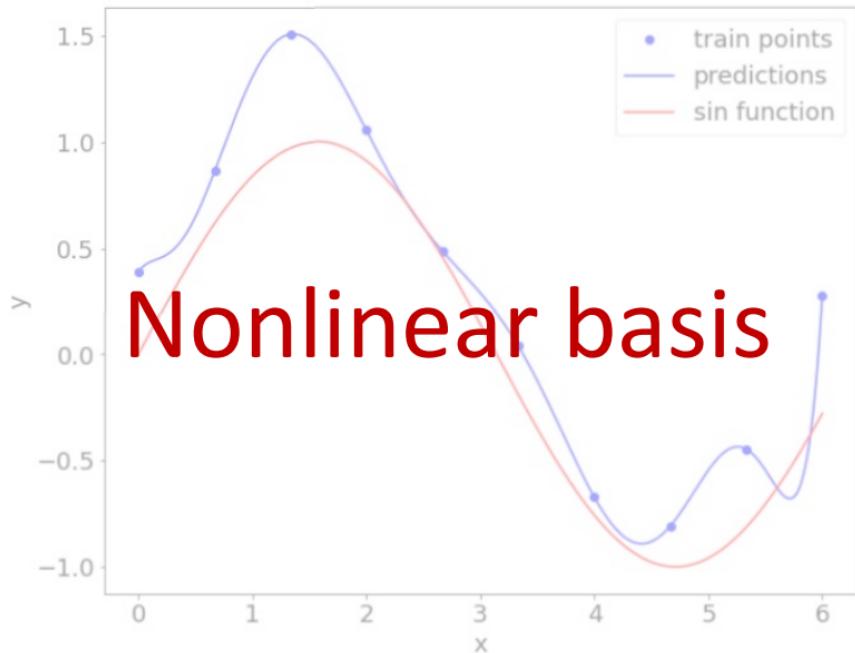
$$|\mathcal{F}| = W^d$$

$\therefore$  generalization gap is at most  $\varepsilon$  with  $n \geq$

*This  $\varepsilon$  maybe  $\varepsilon^2$  without realizability*

$$\frac{\ln(|\mathcal{F}|/\delta)}{\varepsilon} \text{ samples}$$
$$= \frac{d \ln(W/\delta)}{\varepsilon} \text{ samples}$$

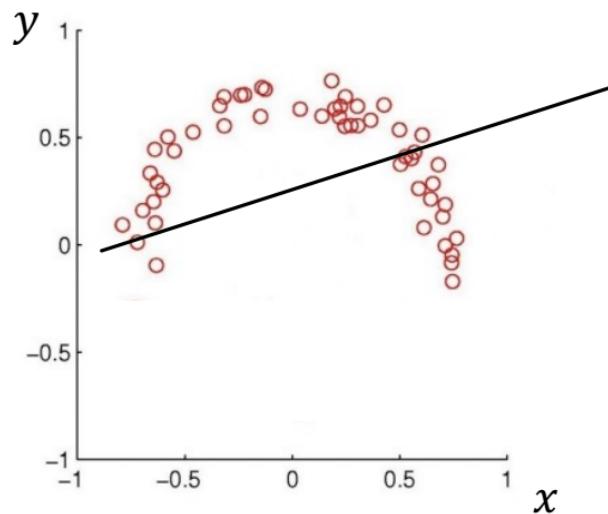
A useful rule of thumb: to guarantee generalization, make sure that your training data set size  $n$  is at least linear in the number  $d$  of free parameters in the function that you're trying to learn.



# What if a linear model is not a good fit?

Let's go back to the regression setup (output  $y \in R$ ).

A linear model could be a bad fit for the following data:



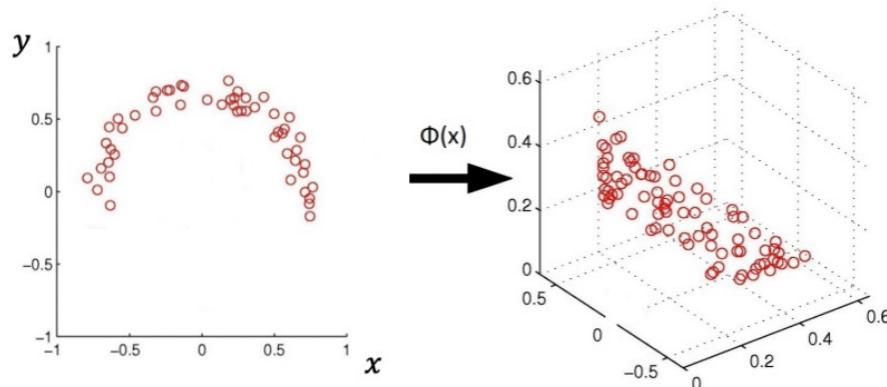
# A solution: nonlinearly transformed features

## 1. Use a nonlinear mapping

$$\phi(x) : x \in \mathbb{R}^d \rightarrow z \in \mathbb{R}^M$$

to transform the data to a more complicated feature space

## 2. Then apply linear regression (hope: linear model is a better fit for the new feature space).



# A solution: nonlinearly transformed features

## 1. Use a nonlinear mapping

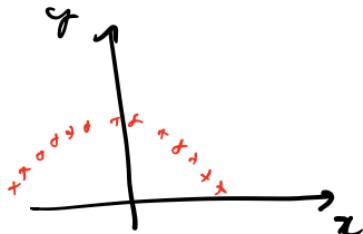
$$\phi(x) : x \in \mathbb{R}^d \rightarrow z \in \mathbb{R}^M$$

to transform the data to a more complicated feature space

## 2. Then apply linear regression (hope: linear model is a better fit for the new feature space).

E.g. consider  $y = 1 - x^2$ ,  $x \in [0, 1]$

$$\phi(x) = \begin{pmatrix} 1 \\ x \\ x^2 \end{pmatrix} \in \mathbb{R}^3$$



$$y = w^T \phi(x)$$

$$\text{for } w = (1, 0, -1)$$

# Regression with nonlinear basis

**Model:**  $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$  where  $\mathbf{w} \in \mathbb{R}^M$

**Objective:**

$$\text{RSS}(\mathbf{w}) = \sum_{i=1}^n (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2$$

**Similar least square solution:**

$$\mathbf{w}^* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \quad \text{where} \quad \Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_n)^T \end{pmatrix} \in \mathbb{R}^{n \times M}$$

# Example

**Polynomial basis functions for  $d = 1$**

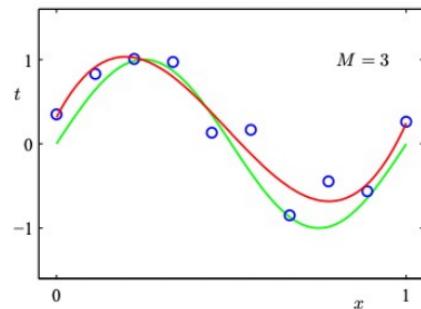
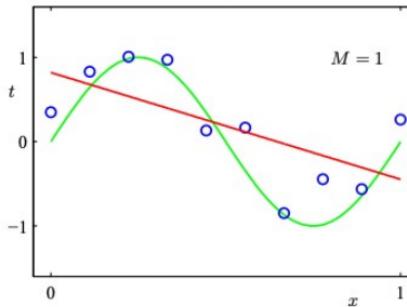
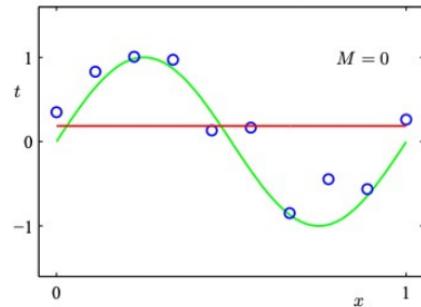
$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow f(x) = w_0 + \sum_{m=1}^M w_m x^m$$

Learning a linear model in the new space

= learning an *M-degree polynomial model* in the original space

# Example

Fitting a noisy sine function with a polynomial ( $M = 0, 1, \text{ or } 3$ ):



See Colab notebook

## Why nonlinear?

Can I use a fancy **linear feature map**?

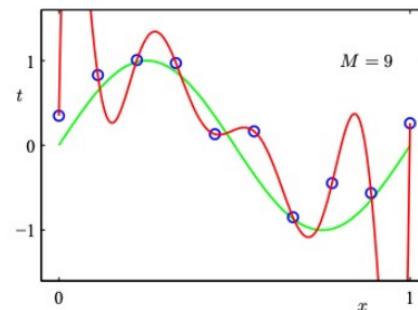
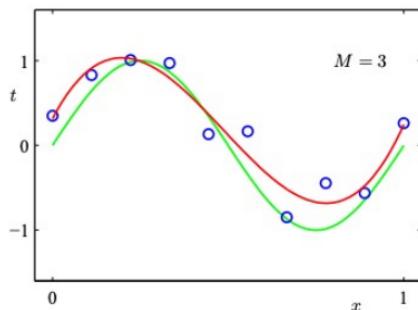
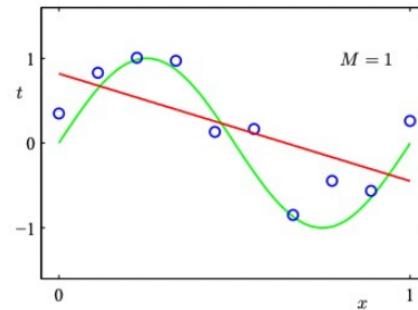
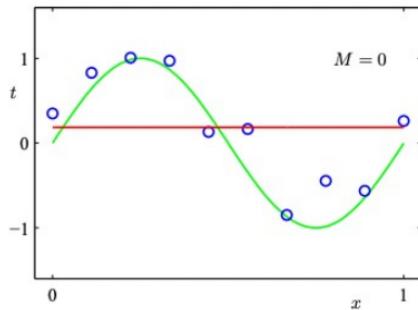
$$\phi(\mathbf{x}) = \begin{bmatrix} x_1 - x_2 \\ 3x_4 - x_3 \\ 2x_1 + x_4 + x_5 \\ \vdots \end{bmatrix} = \mathbf{A}\mathbf{x} \quad \text{for some } \mathbf{A} \in \mathbb{R}^{M \times d}$$



Overfitting and  
Regularization

# Should we use a very complicated mapping?

**Ex: fitting a noisy sine function with a polynomial:**



See Colab notebook

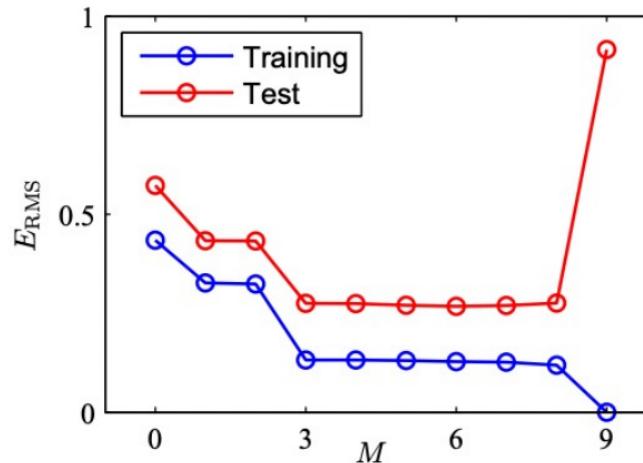
# Underfitting and overfitting

$M \leq 2$  is *underfitting* the data

- large training error
- large test error

$M \geq 9$  is *overfitting* the data

- small training error
- **large test error**

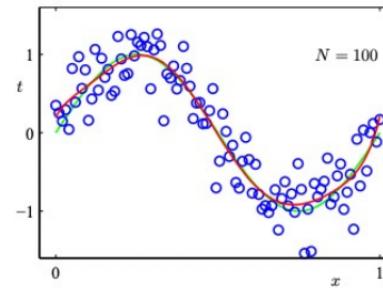
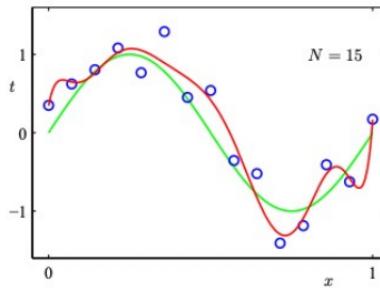
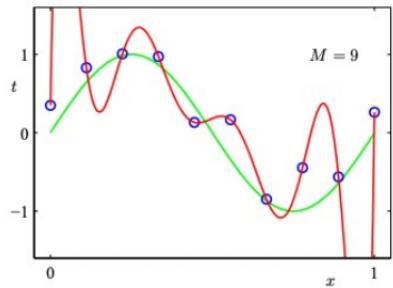


*More complicated models  $\Rightarrow$  larger gap between training and test error*

How to prevent overfitting?

See Colab notebook

# Method 1: More data!!



*More data  $\Rightarrow$  smaller gap between training and test error*

See Colab notebook

## Method 2: Control model complexity

For polynomial basis, the **degree  $M$**  clearly controls the complexity

- use **cross-validation** to pick hyper-parameter  $M$

Cross-validation: Explored in HW1. Idea is to do a three-way split in addition to training set/test set, and tune hyperparameters on a *validation set*.

When  $M$  or in general  $\Phi$  is fixed, are there still other ways to control complexity?

## Magnitude of the weights

Least square solution for the polynomial example:

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
$w_0$	0.19	0.82	0.31	0.35
$w_1$		-1.27	7.99	232.37
$w_2$			-25.43	-5321.83
$w_3$			17.37	48568.31
$w_4$				-231639.30
$w_5$				640042.26
$w_6$				-1061800.52
$w_7$				1042400.18
$w_8$				-557682.99
$w_9$				125201.43

Intuitively, **large weights  $\Rightarrow$  more complex model**

See Colab notebook

# How to make the weights small?

**Regularized linear regression:** new objective

$$G(\mathbf{w}) = \text{RSS}(\mathbf{w}) + \lambda\psi(\mathbf{w})$$

Goal: find  $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} G(\mathbf{w})$

- $\psi : \mathbb{R}^d \rightarrow \mathbb{R}^+$  is the *regularizer*
  - measure how complex the model  $\mathbf{w}$  is, penalize complex models
  - common choices:  $\|\mathbf{w}\|_2^2$ ,  $\|\mathbf{w}\|_1$ , etc.

# How to make the weights small?

**Regularized linear regression:** new objective

$$G(\mathbf{w}) = \text{RSS}(\mathbf{w}) + \lambda\psi(\mathbf{w})$$

Goal: find  $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} G(\mathbf{w})$

- $\psi : \mathbb{R}^d \rightarrow \mathbb{R}^+$  is the *regularizer*
  - measure how complex the model  $\mathbf{w}$  is, penalize complex models
  - common choices:  $\|\mathbf{w}\|_2^2$ ,  $\|\mathbf{w}\|_1$ , etc.
- $\lambda > 0$  is the *regularization coefficient*
  - $\lambda = 0$ , no regularization
  - $\lambda \rightarrow +\infty$ ,  $\mathbf{w} \rightarrow \operatorname{argmin}_{\mathbf{w}} \psi(\mathbf{w})$
  - i.e. control **trade-off** between training error and complexity

# $\ell_2$ regularization with non-linear basis: The effect of $\lambda$

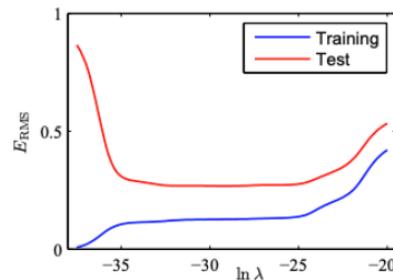
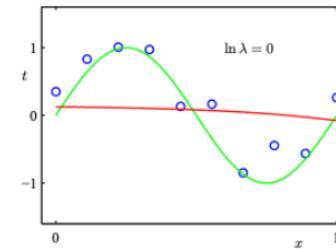
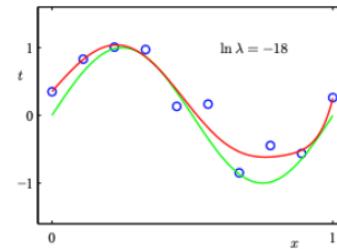
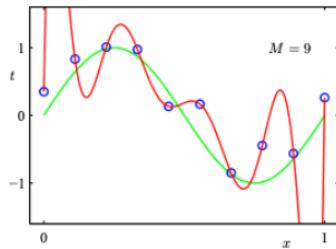
when we increase regularization coefficient  $\lambda$

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
$w_0$	0.35	0.35	0.13
$w_1$	232.37	4.74	-0.05
$w_2$	-5321.83	-0.77	-0.06
$w_3$	48568.31	-31.97	-0.06
$w_4$	-231639.30	-3.89	-0.03
$w_5$	640042.26	55.28	-0.02
$w_6$	-1061800.52	41.32	-0.01
$w_7$	1042400.18	-45.95	-0.00
$w_8$	-557682.99	-91.53	0.00
$w_9$	125201.43	72.68	0.01

See Colab notebook

# $\ell_2$ regularization with non-linear basis : A tradeoff

when we increase regularization coefficient  $\lambda$



See Colab notebook

## Why is regularization useful?

If you don't have sufficient data to fit your more expressive model, then ERM will overfit.  
**Regularization helps with generalization.**

So should it not be useful in many practical settings, where we have enough data?

## Why is regularization useful?

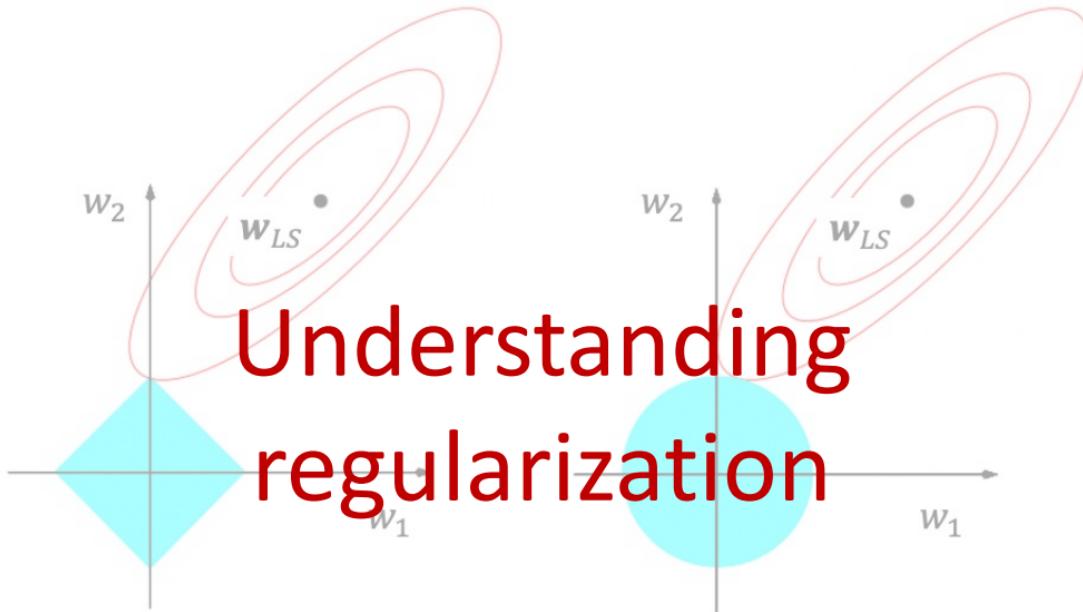
If you don't have sufficient data to fit your more expressive model, then ERM will overfit.  
**Regularization helps with generalization.**

So should it not be useful in many practical settings, where we have enough data?

In general, a viewpoint is that *we should always be trying to fit a more expressive model if possible*. We want our function class to be rich enough that we could almost overfit if we are not careful.

Since we're often in this regime where the models we want to fit are more and more complex, regularization is very useful to help generalization (it's also a relatively simple knob to control).

# Understanding regularization



# How to solve the regularized objective $G(\mathbf{w})$ ?

Let's go back to the original linear model.

**Simple for  $\ell_2$  regularization,**  $\psi(\mathbf{w}) = \|\mathbf{w}\|_2^2$ :

$$G(\mathbf{w}) = \text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

$$\begin{aligned}\nabla G(\mathbf{w}) &= 2(\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y}) + 2\lambda \mathbf{w} = 0 \\ \Rightarrow (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} &= \mathbf{X}^T \mathbf{y} \\ \Rightarrow \mathbf{w}^* &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

Linear regression with  $\ell_2$  regularization is also known as **ridge regression**.

For other regularizers, as long as it's **convex**, standard optimization algorithms can be applied.

## Aside: Least-squares when $X^T X$ is not invertible

When  $X^T X$  is not invertible  $w_{LS} = (X^T X)^{-1} X^T y$  is not defined

This could happen when:

①  $\infty$  many  $w$  s.t.  $Xw = y \rightarrow$  Let's look at this case

② no such  $w$  s.t.  $Xw = y$

① can happen when  $n < d$  (do not have enough data to learn)

What does  $l_2$  regularization do here?

$$h(w) = \underbrace{\|Xw - y\|_2^2}_0 + \lambda \|w\|_2^2$$

for all  $w$  s.t.  $Xw = y$

$\therefore l_2$  regularization chooses  $w$  with smallest  $\|w\|_2$  s.t.  $Xw = y$ .

## Aside: Least-squares when $\mathbf{X}^T \mathbf{X}$ is not invertible

**Intuition:** what does inverting  $\mathbf{X}^T \mathbf{X}$  do?

**eigendecomposition:**  $\mathbf{X}^T \mathbf{X} = \mathbf{U}^T \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & \lambda_d & 0 \\ 0 & \cdots & 0 & \lambda_{d+1} \end{bmatrix} \mathbf{U}$

where  $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_{d+1} \geq 0$  are **eigenvalues**.

**inverse:**  $(\mathbf{X}^T \mathbf{X})^{-1} = \mathbf{U}^T \begin{bmatrix} \frac{1}{\lambda_1} & 0 & \cdots & 0 \\ 0 & \frac{1}{\lambda_2} & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & \frac{1}{\lambda_d} & 0 \\ 0 & \cdots & 0 & \frac{1}{\lambda_{d+1}} \end{bmatrix} \mathbf{U}$

i.e. just invert the eigenvalues

## Aside: Least-squares when $\mathbf{X}^T \mathbf{X}$ is not invertible

Non-invertible  $\Rightarrow$  some eigenvalues are 0.

**One natural fix: add something positive**

$$\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} = \mathbf{U}^T \begin{bmatrix} \lambda_1 + \lambda & 0 & \cdots & 0 \\ 0 & \lambda_2 + \lambda & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & \lambda_d + \lambda & 0 \\ 0 & \cdots & 0 & \lambda_{d+1} + \lambda \end{bmatrix} \mathbf{U}$$

where  $\lambda > 0$  and  $\mathbf{I}$  is the identity matrix. Now it is invertible:

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} = \mathbf{U}^T \begin{bmatrix} \frac{1}{\lambda_1+\lambda} & 0 & \cdots & 0 \\ 0 & \frac{1}{\lambda_2+\lambda} & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & \frac{1}{\lambda_d+\lambda} & 0 \\ 0 & \cdots & 0 & \frac{1}{\lambda_{d+1}+\lambda} \end{bmatrix} \mathbf{U}$$

## A “Bayesian view” of $\ell_2$ regularization

**Maximum a posteriori probability (MAP) estimation:** A Bayesian generalization of maximum likelihood estimation (MLE).

Let's continue with the linear model, and Q3 from the practice problems for today.

Have training set  $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^d \times \mathbb{R}$

$$y_i = w_*^\top x_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2)$$

$$\therefore \log(p_{\theta}(y|x_i; w_*, \sigma)) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y - w_*^\top x_i)^2}{2\sigma^2}\right)$$

MLE : find  $w$  which maximizes likelihood (i.e.  $p_{\theta}(y_i|x_i; w, \sigma)$ )

## A “Bayesian view” of $\ell_2$ regularization

**Maximum a posteriori probability (MAP) estimation:** A Bayesian generalization of maximum likelihood estimation (MLE).

Let's continue with the linear model, and Q3 from the practice problems for today.

$$\log \left( p(y_i | x_i; w) \right) = -\frac{(y_i - w^\top x_i)^2}{2\sigma^2} + \text{const}$$

The solution is  $w_* = (x^\top x)^{-1} x^\top y$

## A “Bayesian view” of $\ell_2$ regularization

**Maximum a posteriori probability (MAP) estimation:** A Bayesian generalization of maximum likelihood estimation (MLE).

Bayesian view: A **prior** over  $w$

Suppose our prior for  $w$  is  $N(0, \gamma^2 I)$

Now we find the model which maximizes posterior  
probability (MAP)

Posterior  $\propto$  Prior. Likelihood

## A “Bayesian view” of $\ell_2$ regularization

**Maximum a posteriori probability (MAP) estimation:** A Bayesian generalization of maximum likelihood estimation (MLE).

Bayesian view: A **prior** over  $w$

$$\text{Posterior}(w) \propto \underbrace{\prod_{j=1}^d \exp\left(-\frac{w_j^2}{2\gamma^2}\right)}_{w \sim N(0, \gamma^2 I)} \cdot \prod_{i=1}^n \exp\left(-\frac{(y_i - w^\top x_i)^2}{2\sigma^2}\right)$$

$$\log(\text{Posterior}(w)) = -\frac{\|w\|^2}{2\gamma^2} - \sum_{i=1}^n (y_i - w^\top x_i)^2$$

max  $(\log(\text{Posterior}))$  is same as  $\min_b b(w)$  for  $\Psi(w) = \|w\|^2$

## An equivalent form, and a “Frequentist view”

“Frequentist” approach to justifying regularization is to argue that if the true model has a specific property, then regularization will allow you to recover a good approximation to the true model. In this view, we can equivalently formulate regularization as:

$$\operatorname{argmin}_{\mathbf{w}} \text{RSS}(\mathbf{w}) \quad \text{subject to } \psi(\mathbf{w}) \leq \beta$$

(consider  $\psi(\mathbf{w}) = \|\mathbf{w}\|_2^2$ )

where  $\beta$  is some hyper-parameter.

Finding the solution becomes a *constrained optimization problem*.

Choosing either  $\lambda$  or  $\beta$  can be done by cross-validation.

## Encouraging sparsity: $\ell_0$ regularization

Continuing from the frequentist view, having small norm is one possible structure to impose on the model. Another very common one is **sparsity**.

**Sparsity of  $w$ :** Number of non-zero coefficients in  $w$ . Same as  $\|w\|_0$

E.g.  $w = [1, 0, -1, 0, 0.2, 0, 0]$  is 3-sparse

$\|w\|_0$  is not a norm.

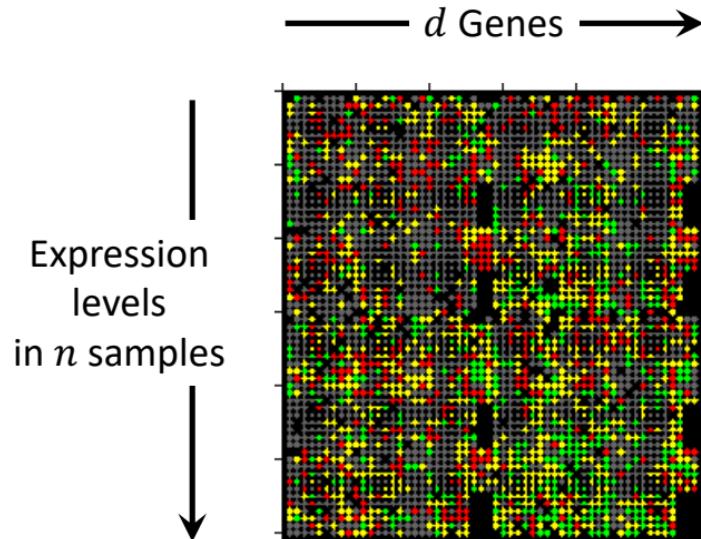
But behaves like  $\lim_{p \rightarrow 0} \|w\|_p$

# Encouraging sparsity: $\ell_0$ regularization

**Sparsity of  $w$ :** Number of non-zero coefficients in  $w$ . Same as  $\|w\|_0$

Advantage:

- Sparse models are a natural inductive bias in many settings. In many applications we have numerous possible features, only some of which may have any relationship with the label.



Suppose we want to fit a linear model from gene expression to an outcome (disease, phenotype etc.).

$d$  is huge, but likely that only a few genes are related.

## Encouraging sparsity: $\ell_0$ regularization

**Sparsity of  $w$ :** Number of non-zero coefficients in  $w$ . Same as  $\|w\|_0$

Advantage:

- Sparse models are a natural inductive bias in many settings. In many applications we have numerous possible features, only some of which may have any relationship with the label.
- Sparse models may also be more **interpretable**. They could narrow down a small number of features which carry a lot of signal.

E.g.      $w = [1.5, 0, -1.1, 0, 0.25, 0, 0]$  is more interpretable than,  
             $w = [1, 0.2, -1.3, 0.15, 0.2, 0.05, 0.12]$

For a sparse model, it could be easier to understand the model. It is also easier to verify whether the features which have a high weight have a relation with the outcome (they are not spurious artifacts of the data).

# Encouraging sparsity: $\ell_0$ regularization

**Sparsity of  $w$ :** Number of non-zero coefficients in  $w$ . Same as  $\|w\|_0$

Advantage:

- Sparse models are a natural inductive bias in many settings. In many applications we have numerous possible features, only some of which may have any relationship with the label.
- Sparse models may also be more **interpretable**. They could narrow down a small number of features which carry a lot of signal.
- Data required to learn sparse model maybe significantly less than to learn dense model.

We'll see more on the third point next.

# $\ell_0$ regularization: The good, the bad and the ugly

# $\ell_0$ regularization: The good, the bad and the ugly

# $\ell_0$ regularization: The good, the bad and the ugly

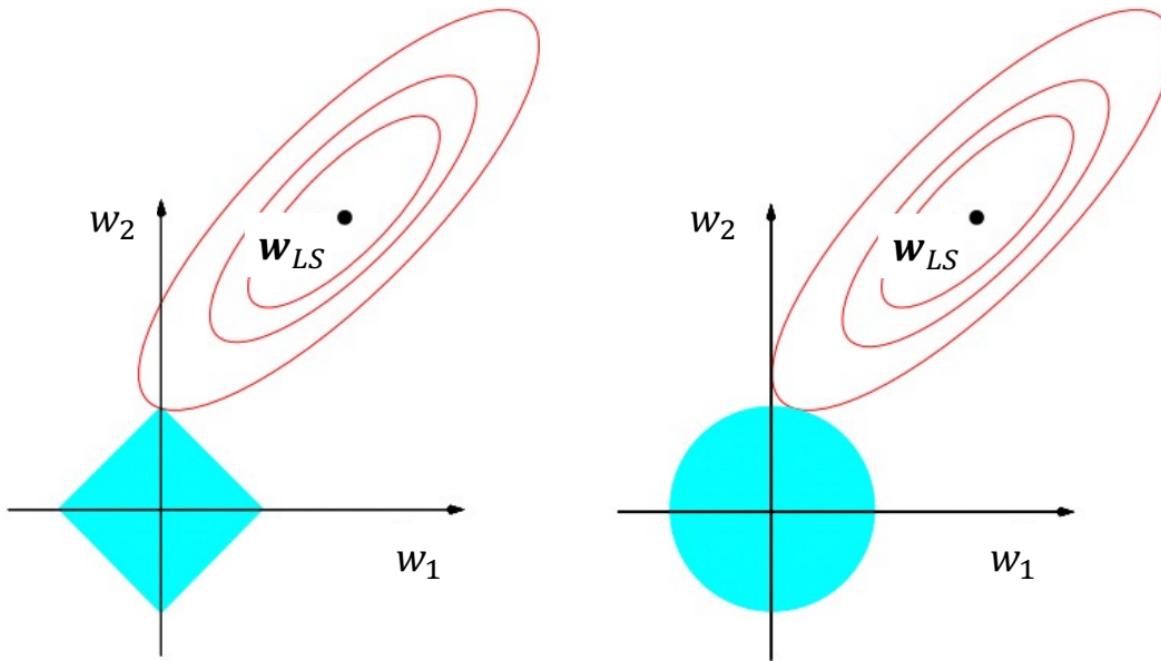
# $\ell_0$ regularization: The good, the bad and the ugly

$\ell_1$  regularization as a proxy for  $\ell_0$  regularization

$\ell_1$  regularization as a proxy for  $\ell_0$  regularization

# Why does $\ell_1$ regularization encourage sparse solutions?

Optimization problem:  $\operatorname{argmin}_w \text{RSS}(w)$ , subject to  $\psi(w) \leq \beta$

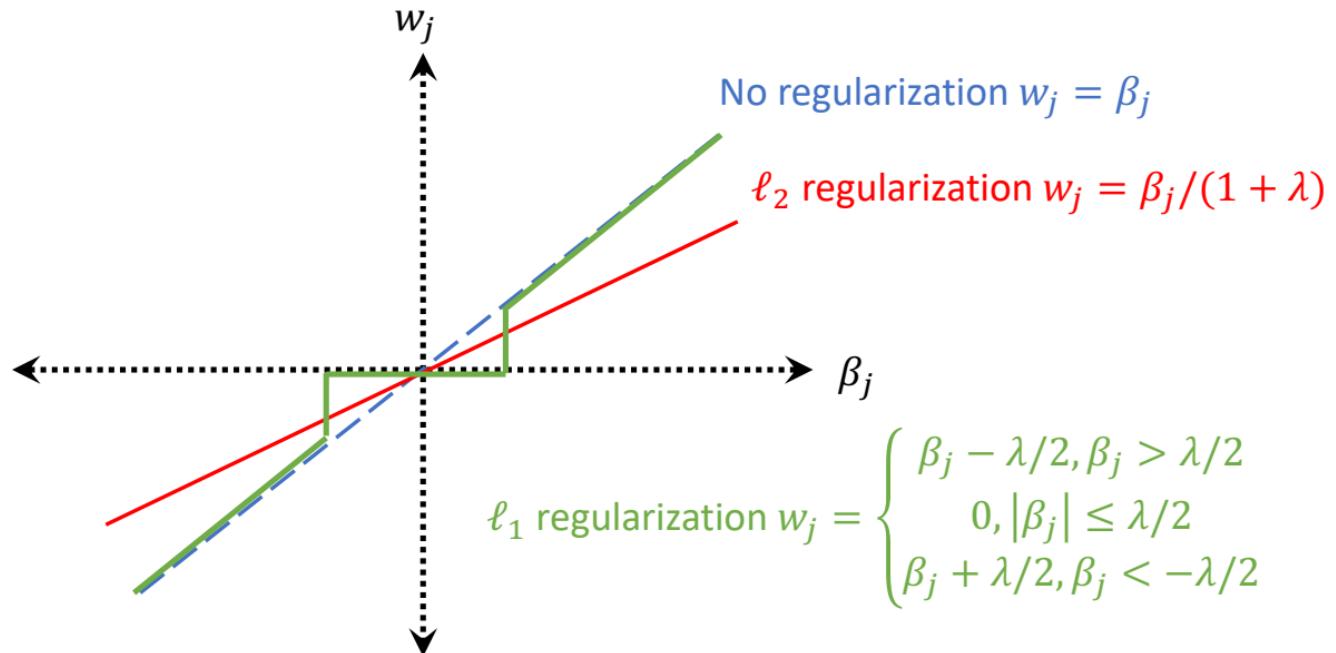


$\ell_1$  and  $\ell_2$  regularization for the “isotropic” case

# $\ell_1$ and $\ell_2$ regularization for the “isotropic” case

Summary: Isotropic case ( $\mathbf{X}^T \mathbf{X} = \mathbf{I}$ ).

Let  $\beta_j = \mathbf{X}_{(i)}^T \mathbf{y}$



## Implicit regularization

So far, we explicitly added a  $\psi(\mathbf{w})$  term to our objective function to regularize.

In many cases, the optimization algorithm we use can themselves act as regularizers, favoring some solutions over others.

Currently a very active area of research, you'll see more in the homework.



# Multiclass classification

## Setup

- input (feature vector):  $\boldsymbol{x} \in \mathbb{R}^d$
- output (label):  $y \in [C] = \{1, 2, \dots, C\}$
- goal: learn a mapping  $f : \mathbb{R}^d \rightarrow [C]$

## Examples:

- recognizing digits ( $C = 10$ ) or letters ( $C = 26$  or  $52$ )
- predicting weather: sunny, cloudy, rainy, etc
- predicting image category: ImageNet dataset ( $C \approx 20K$ )

## Linear models: Binary to multiclass

Step 1: *What should a linear model look like for multiclass tasks?*

Note: a linear model for binary tasks (switching from  $\{-1, +1\}$  to  $\{1, 2\}$ )

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} \geq 0 \\ 2 & \text{if } \mathbf{w}^T \mathbf{x} < 0 \end{cases}$$

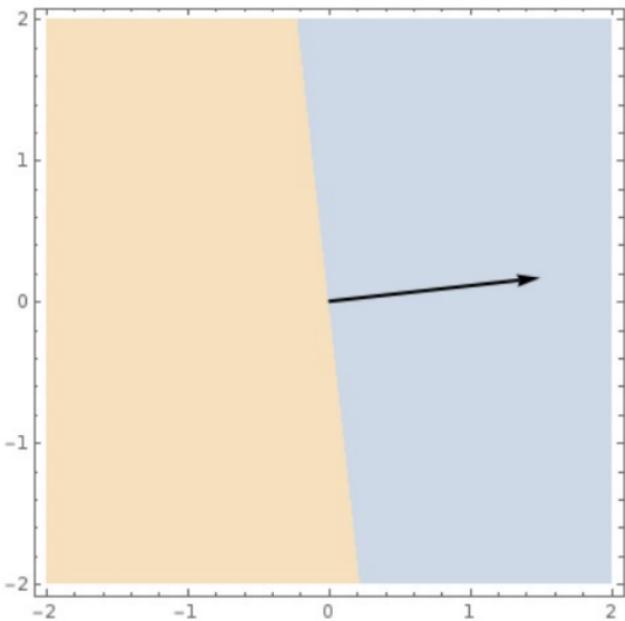
can be written as

$$\begin{aligned} f(\mathbf{x}) &= \begin{cases} 1 & \text{if } \mathbf{w}_1^T \mathbf{x} \geq \mathbf{w}_2^T \mathbf{x} \\ 2 & \text{if } \mathbf{w}_2^T \mathbf{x} > \mathbf{w}_1^T \mathbf{x} \end{cases} \\ &= \operatorname{argmax}_{k \in \{1, 2\}} \mathbf{w}_k^T \mathbf{x} \end{aligned}$$

for any  $\mathbf{w}_1, \mathbf{w}_2$  s.t.  $\mathbf{w} = \mathbf{w}_1 - \mathbf{w}_2$

Think of  $\mathbf{w}_k^T \mathbf{x}$  as a **score for class  $k$** .

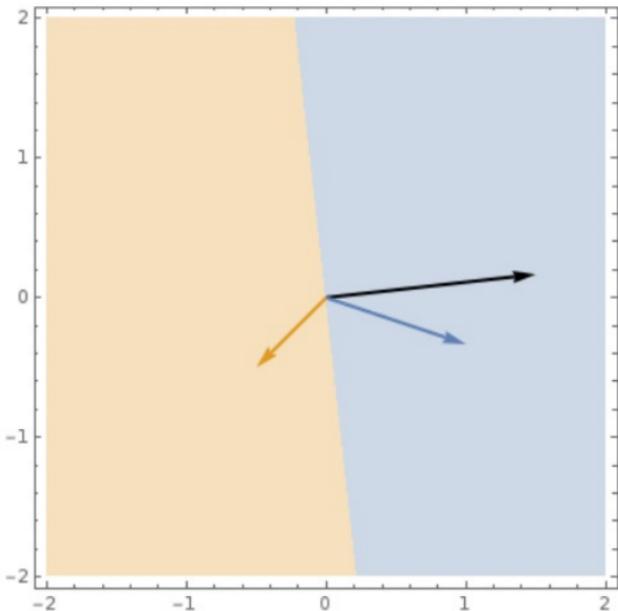
## Linear models: Binary to multiclass



$$\mathbf{w} = \left(\frac{3}{2}, \frac{1}{6}\right)$$

- Blue class:  
 $\{\mathbf{x} : \mathbf{w}^T \mathbf{x} \geq 0\}$
- Orange class:  
 $\{\mathbf{x} : \mathbf{w}^T \mathbf{x} < 0\}$

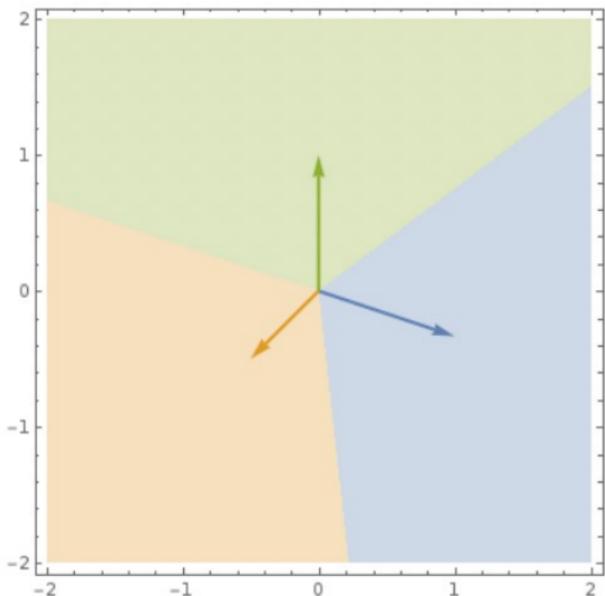
## Linear models: Binary to multiclass



$$\begin{aligned}\mathbf{w} &= \left(\frac{3}{2}, \frac{1}{6}\right) = \mathbf{w}_1 - \mathbf{w}_2 \\ \mathbf{w}_1 &= \left(1, -\frac{1}{3}\right) \\ \mathbf{w}_2 &= \left(-\frac{1}{2}, -\frac{1}{2}\right)\end{aligned}$$

- Blue class:  
 $\{\mathbf{x} : 1 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$
- Orange class:  
 $\{\mathbf{x} : 2 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$

## Linear models: Binary to multiclass



$$\begin{aligned}w_1 &= (1, -\frac{1}{3}) \\w_2 &= (-\frac{1}{2}, -\frac{1}{2}) \\w_3 &= (0, 1)\end{aligned}$$

- **Blue class:**  
 $\{x : 1 = \text{argmax}_k w_k^T x\}$
- **Orange class:**  
 $\{x : 2 = \text{argmax}_k w_k^T x\}$
- **Green class:**  
 $\{x : 3 = \text{argmax}_k w_k^T x\}$

## Linear models: Function class

$$\begin{aligned}\mathcal{F} &= \left\{ f(\mathbf{x}) = \operatorname{argmax}_{k \in [C]} \mathbf{w}_k^T \mathbf{x} \mid \mathbf{w}_1, \dots, \mathbf{w}_C \in \mathbb{R}^d \right\} \\ &= \left\{ f(\mathbf{x}) = \operatorname{argmax}_{k \in [C]} (\mathbf{W} \mathbf{x})_k \mid \mathbf{W} \in \mathbb{R}^{C \times d} \right\}\end{aligned}$$

Next, let's try to generalize the loss functions. Focus on the logistic loss today.

## Multinomial logistic regression: a probabilistic view

Observe: for binary logistic regression, with  $\mathbf{w} = \mathbf{w}_1 - \mathbf{w}_2$ :

$$\mathbb{P}(y = 1 \mid \mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = \frac{e^{\mathbf{w}_1^T \mathbf{x}}}{e^{\mathbf{w}_1^T \mathbf{x}} + e^{\mathbf{w}_2^T \mathbf{x}}} \propto e^{\mathbf{w}_1^T \mathbf{x}}$$

Naturally, for multiclass:

$$\mathbb{P}(y = k \mid \mathbf{x}; \mathbf{W}) = \frac{e^{\mathbf{w}_k^T \mathbf{x}}}{\sum_{k' \in [C]} e^{\mathbf{w}_{k'}^T \mathbf{x}}} \propto e^{\mathbf{w}_k^T \mathbf{x}}$$

This is called the *softmax function*.

## Let's find the **MLE**

Maximize probability of seeing labels  $y_1, \dots, y_n$  given  $\mathbf{x}_1, \dots, \mathbf{x}_n$

$$P(\mathbf{W}) = \prod_{i=1}^n \mathbb{P}(y_i \mid \mathbf{x}_i; \mathbf{W}) = \prod_{i=1}^n \frac{e^{\mathbf{w}_{y_i}^T \mathbf{x}_i}}{\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}_i}}$$

By taking **negative log**, this is equivalent to minimizing

$$F(\mathbf{W}) = \sum_{i=1}^n \ln \left( \frac{\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}_i}}{e^{\mathbf{w}_{y_i}^T \mathbf{x}_i}} \right) = \sum_{i=1}^n \ln \left( 1 + \sum_{k \neq y_i} e^{(\mathbf{w}_k - \mathbf{w}_{y_i})^T \mathbf{x}_i} \right)$$

This is the *multiclass logistic loss*.

When  $C = 2$ , this is the same as binary logistic loss.

Let's find the MLE

## Next, optimization

Apply **SGD**: what is the gradient of

$$F(\mathbf{W}) = \ln \left( 1 + \sum_{k \neq y_i} e^{(\mathbf{w}_k - \mathbf{w}_{y_i})^T \mathbf{x}_i} \right) ?$$

It's a  $C \times d$  matrix. Let's focus on the  $k$ -th row:

If  $k \neq y_i$ :

$$\nabla_{\mathbf{w}_k^T} F(\mathbf{W}) = \frac{e^{(\mathbf{w}_k - \mathbf{w}_{y_i})^T \mathbf{x}_i}}{1 + \sum_{k \neq y_i} e^{(\mathbf{w}_k - \mathbf{w}_{y_i})^T \mathbf{x}_i}} \mathbf{x}_i^T = \mathbb{P}(k \mid \mathbf{x}_i; \mathbf{W}) \mathbf{x}_i^T$$

else:

$$\nabla_{\mathbf{w}_k^T} F(\mathbf{W}) = \frac{- \left( \sum_{k \neq y_i} e^{(\mathbf{w}_k - \mathbf{w}_{y_i})^T \mathbf{x}_i} \right)}{1 + \sum_{k \neq y_i} e^{(\mathbf{w}_k - \mathbf{w}_{y_i})^T \mathbf{x}_i}} \mathbf{x}_i^T = (\mathbb{P}(y_i \mid \mathbf{x}_i; \mathbf{W}) - 1) \mathbf{x}_i^T$$

# SGD for multinomial logistic regression

Initialize  $\mathbf{W} = \mathbf{0}$  (or randomly). Repeat:

- ① pick  $i \in [n]$  uniformly at random
- ② update the parameters

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \begin{pmatrix} \mathbb{P}(y = 1 \mid \mathbf{x}_i; \mathbf{W}) \\ \vdots \\ \mathbb{P}(y = y_i \mid \mathbf{x}_i; \mathbf{W}) - 1 \\ \vdots \\ \mathbb{P}(y = C \mid \mathbf{x}_i; \mathbf{W}) \end{pmatrix} \mathbf{x}_i^T$$

Think about why the algorithm makes sense intuitively.

## Probabilities -> Prediction

Having learned  $\mathbf{W}$ , we can either

- make a *deterministic* prediction  $\operatorname{argmax}_{k \in [C]} \mathbf{w}_k^T \mathbf{x}$
- make a *randomized* prediction according to  $\mathbb{P}(k \mid \mathbf{x}; \mathbf{W}) \propto e^{\mathbf{w}_k^T \mathbf{x}}$

In either case, **(expected) mistake is bounded by logistic loss**

- deterministic

$$\mathbb{E}[f(\mathbf{x}) \neq y] \leq \log_2 \left( 1 + \sum_{k \neq y} e^{(\mathbf{w}_k - \mathbf{w}_y)^T \mathbf{x}} \right)$$

$\log_2(1 + e^x) \geq 1$   
for  $x \geq 0$

- randomized

$$\mathbb{E}[\mathbb{I}[f(\mathbf{x}) \neq y]] = 1 - \mathbb{P}(y \mid \mathbf{x}; \mathbf{W}) \leq -\ln \mathbb{P}(y \mid \mathbf{x}; \mathbf{W})$$