

CSCI 567: Machine Learning

Vatsal Sharan
Spring 2026

Lecture 3, January 30

HW1 & Logistics

- HW1 due in less than 2 weeks (2/11 at midnight)
- Max 1 late day per HW
- Working in groups:
 - The point of working in groups is not to split problems among 2 people
 - You should both work **together** on the questions, e.g.
 - Figure out math together for theory questions
 - Pair programming for programming questions
- Next week (2/6): Review of concepts (calculus for matrices and vector), problem solving, pytorch

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**.*

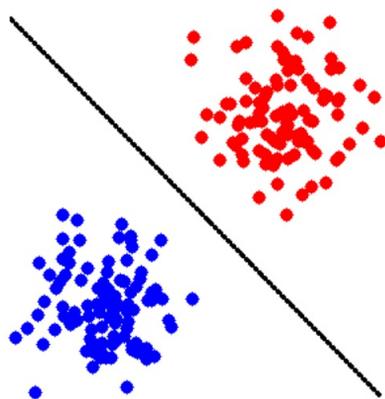
Summary: **Optimization** methods

- **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

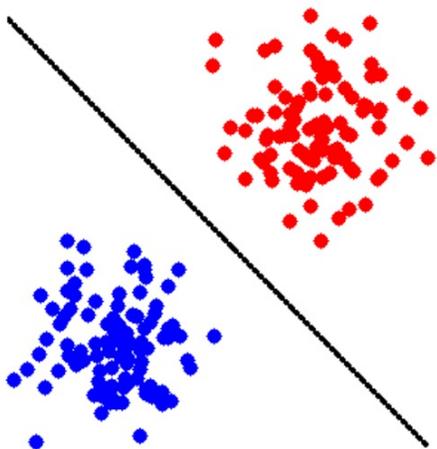
Binary classification:

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



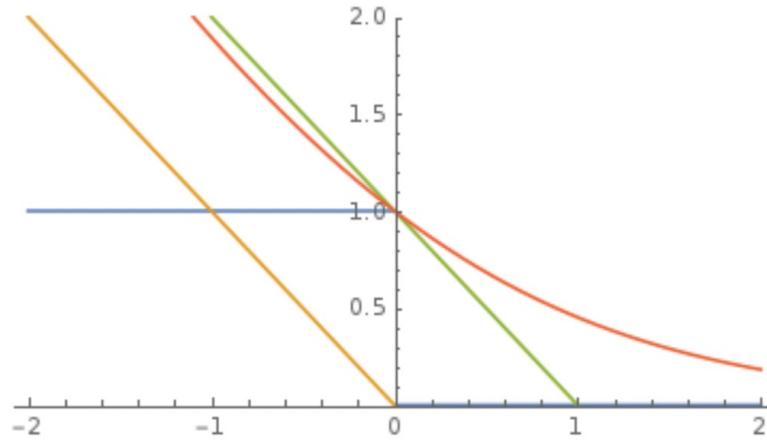
Representation

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:

$$\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(y_i \mathbf{w}^\top \mathbf{x}_i)$$

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_n given x_1, \dots, x_n , as a function of some w ?

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

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 . (*independent & identically distributed*)
- 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 2 on HW1.

Assumptions for today's theory

Finite sized function class

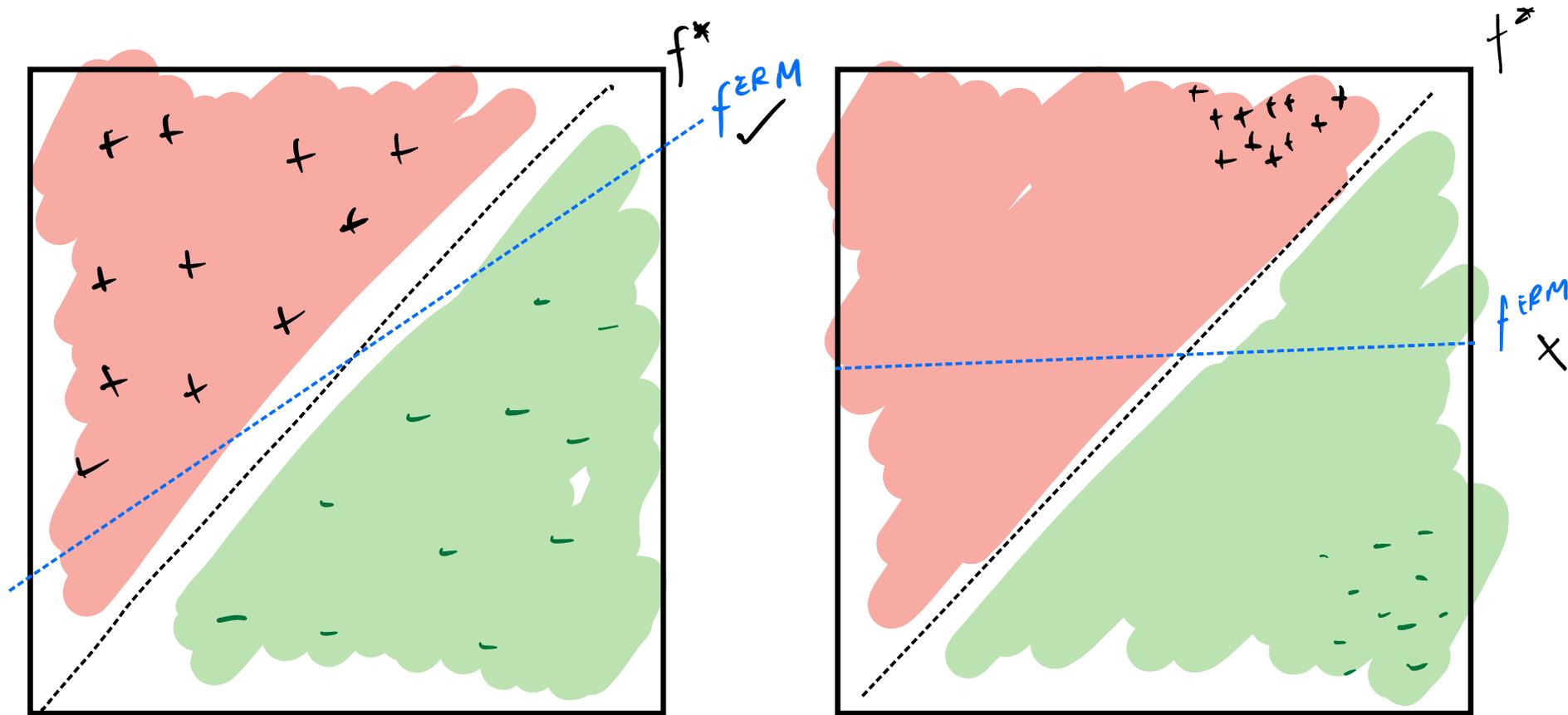
Def: A function class \tilde{F} is finite-sized if $|\tilde{F}|$ is finite.

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

$$|\tilde{F}| = 3^d$$

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

Intuition: When does ERM generalize?



Distribution over $x = (x_1, x_2)$
is uniform

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

$$f_S^{\text{ERM}} = \operatorname{argmin}_{f \in F} \left(\frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) \right).$$

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

Ex. if $\epsilon = 0.1$, $\delta = 0.1$, then n with $n \geq 10 \ln(10/|\mathcal{F}|)$,
with probability 0.9 , $R(f_s^{\text{ERM}}) \leq 0.1$.

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

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

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

Consider some $f' \in \mathcal{F}_{\text{bad}}$

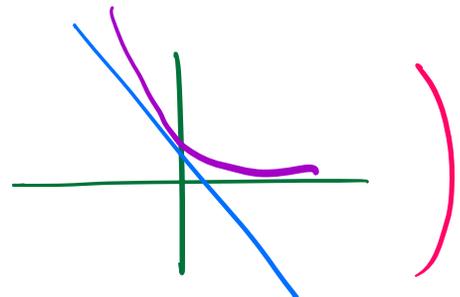
$$\Pr_{S \sim D^n} [f' \text{ is an ERM}] \quad (f^* \text{ gets } \hat{R}_S(f^*) = 0)$$

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

$$= \prod_{i=1}^n \Pr_{x_i \sim D} [f'(x_i) = f^*(x_i)] \quad (\text{i.i.d. assumption, } \Pr(\bar{E}_1 \cap \bar{E}_2) = \Pr(\bar{E}_1) \Pr(\bar{E}_2) \text{ for independent } E_1, E_2)$$

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

$$\leq \prod_{i=1}^n e^{-\epsilon} = e^{-\epsilon n} \quad (\text{since } 1 - x \leq e^{-x})$$



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

Union bound.

$$P_{\mathcal{R}_{S \sim D^n}} \left[\bigcup_{f \in \overline{F}_{\text{bad}}} \{f \text{ is an ERM}\} \right]$$

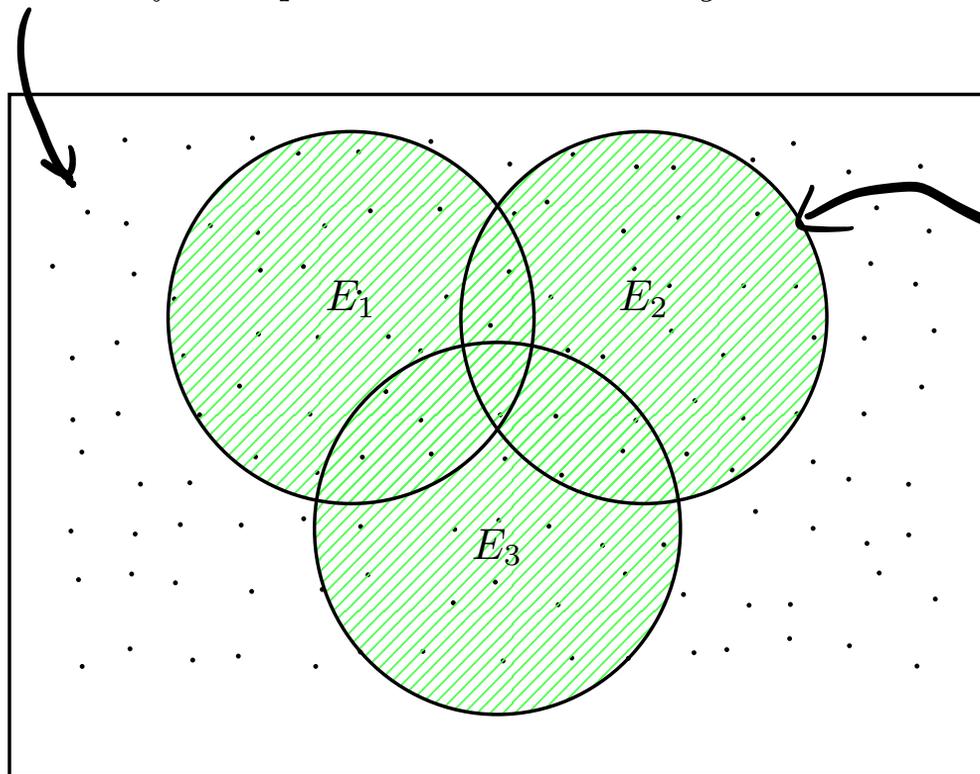
$$\leq \sum_{f \in \overline{F}_{\text{bad}}} P_{\mathcal{R}} [f \text{ is an ERM}]$$

$$\leq \sum_{f \in \overline{F}_{\text{bad}}} e^{-2n} = |\overline{F}_{\text{bad}}| e^{-2n} \leq |F| e^{-\epsilon n}$$

Union bound:

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

Think of each point here as a training set S .



E_2 is the set of all training sets for which some bad hypothesis $f_2 \in \mathcal{H}_{\text{bad}}$ is picked

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

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

$$\therefore \text{if } n \geq \frac{\ln(|F|/\delta)}{\varepsilon}, \text{ w.p. } 1-\delta \quad f_S^{\text{ERM}} \notin F_{\text{bad}}$$

$$\Rightarrow R(f_S^{\text{ERM}}) \leq \varepsilon. \quad \blacksquare$$

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 ϵ^2 instead of ϵ in the denominator). This is called agnostic learning.

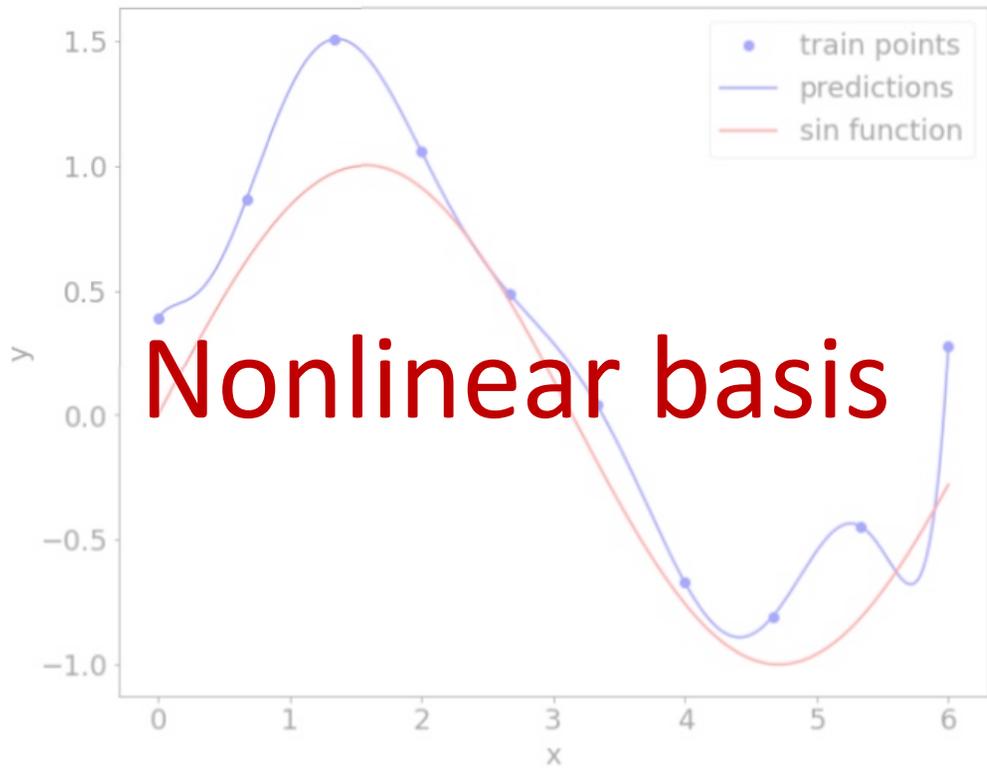
Rule of thumb for generalization

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 each take 3 possible values $\{-1, 0, +1\}$. How much data do we need to have small generalization gap?

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

$$\begin{aligned} \therefore \text{generalization gap of at most } \varepsilon \text{ with } n &\geq \frac{\ln(|\mathcal{F}|/\delta)}{\varepsilon} \text{ samples} \\ &= \frac{d \ln(3/\delta)}{\varepsilon} \text{ samples} \end{aligned}$$

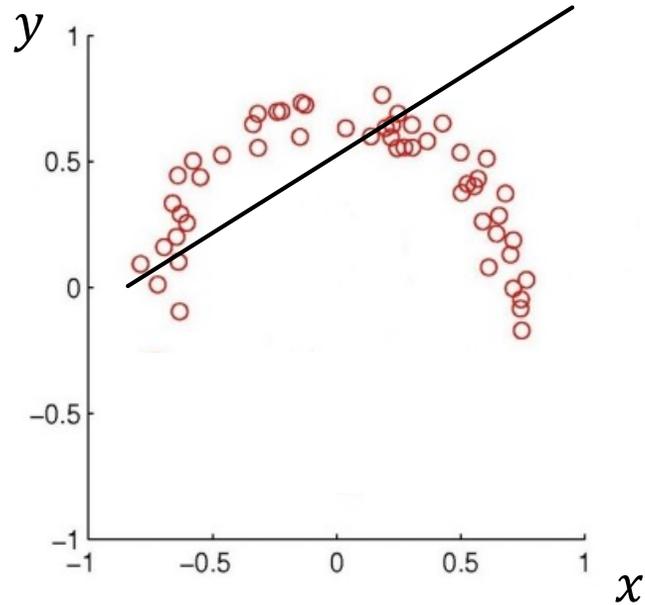
A useful rule of thumb: to guarantee generalization, make sure that your training data set size n is at least linear in the number d of free parameters in the function that you're trying to learn.



What if a linear model is not a good fit?

Let's go back to the regression setup (output $\mathcal{Y} \in R$).

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



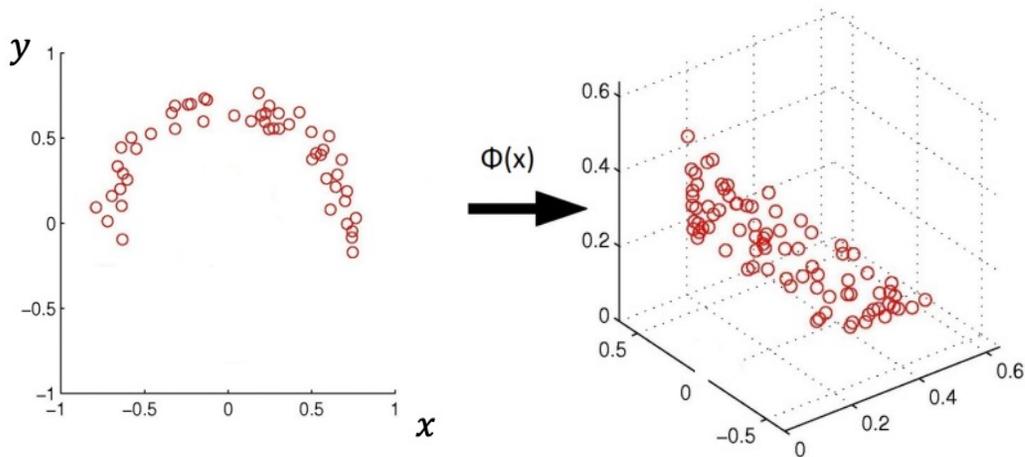
A solution: nonlinearly transformed features

1. Use a nonlinear mapping

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

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

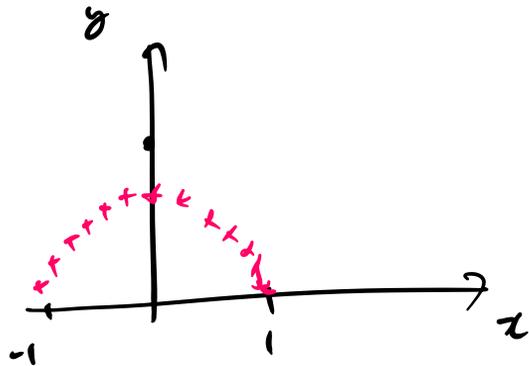
1. Use a nonlinear mapping

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

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 $y = 1 - x^2$, $x \in [-1, 1]$



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

$$y = \mathbf{w}^T \phi(x)$$

$$\text{for } \mathbf{w} = [1, 0, -1]$$

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 $d = 1$

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

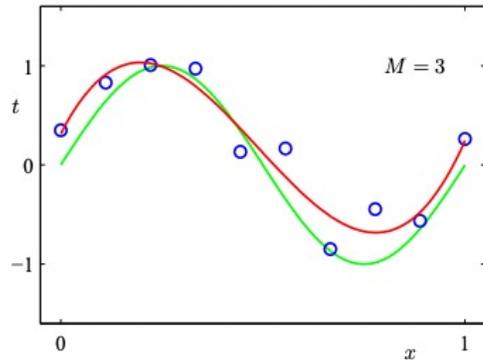
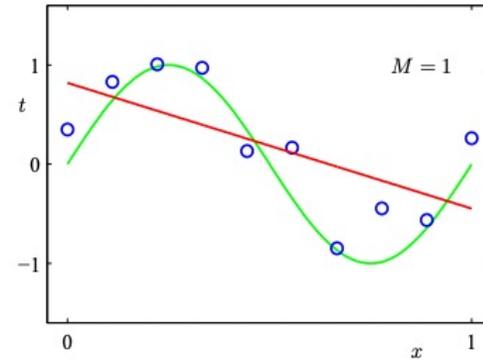
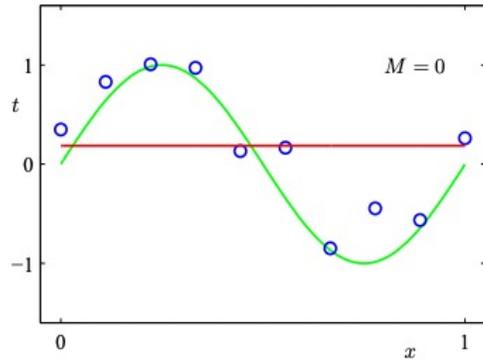
$w^T \phi(x)$

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 ($M = 0, 1, \text{ or } 3$):



Why nonlinear?

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

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

Why nonlinear?

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

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

No, it basically *does nothing* since

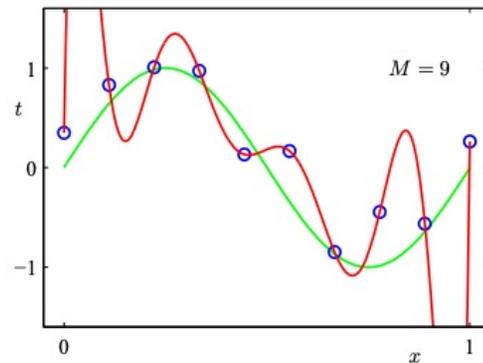
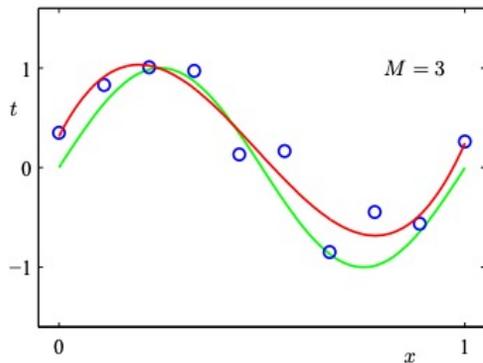
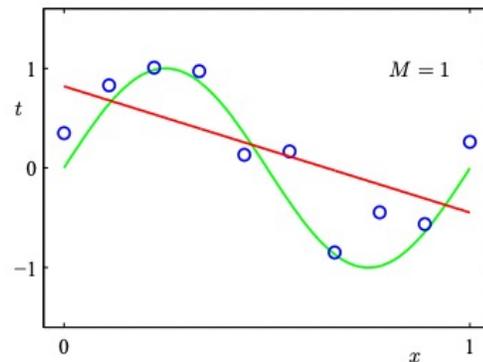
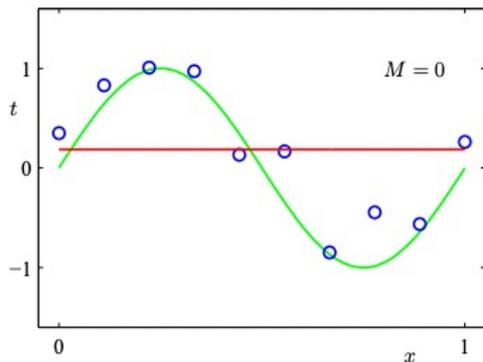
$$\min_{\mathbf{w} \in \mathbb{R}^M} \sum_i (\mathbf{w}^T \mathbf{A}\mathbf{x}_i - y_i)^2 = \min_{\mathbf{w}' \in \text{Im}(\mathbf{A}^T) \subset \mathbb{R}^d} \sum_i (\mathbf{w}'^T \mathbf{x}_i - y_i)^2$$



**Overfitting and
Regularization**

Should we use a very complicated mapping?

Ex: fitting a noisy sine function with a polynomial:



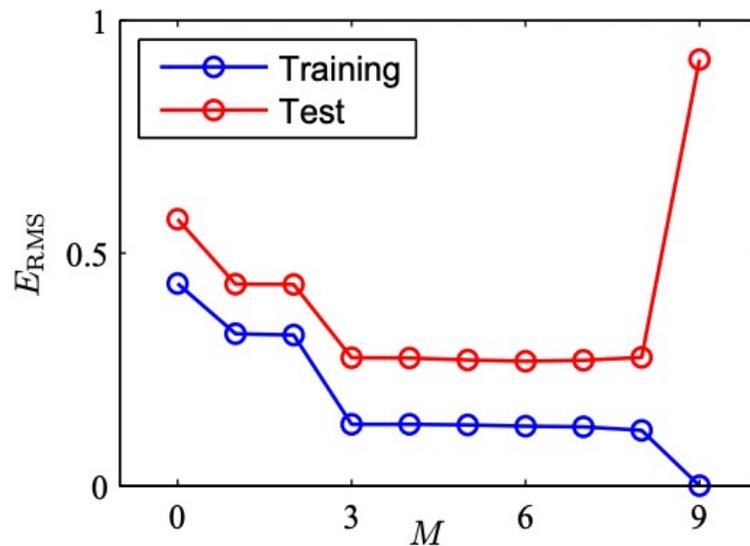
Underfitting and overfitting

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

- large training error
- large test error

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

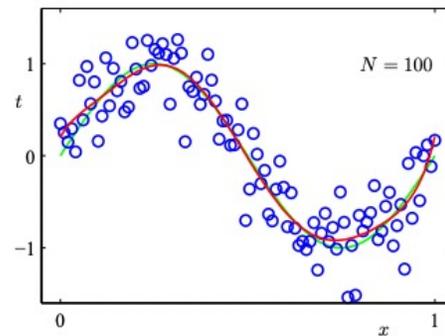
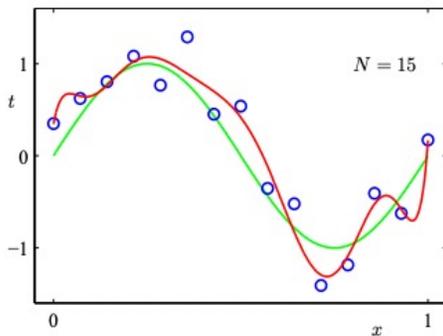
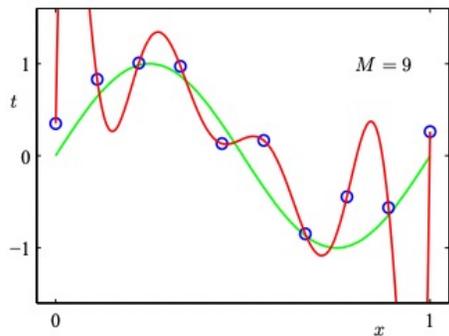
- small training error
- **large test error**



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

How to prevent overfitting?

Method 1: More data!!



More data \Rightarrow smaller gap between training and test error

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 Φ 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**

How to make the weights small?

Regularized linear regression: new objective

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

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

- $\psi : \mathbb{R}^d \rightarrow \mathbb{R}^+$ 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

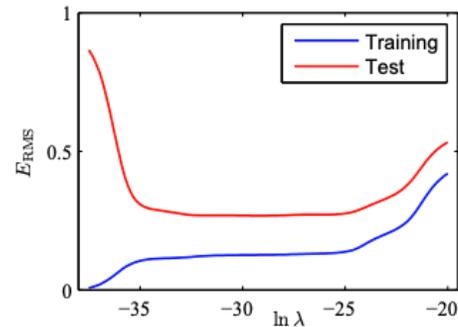
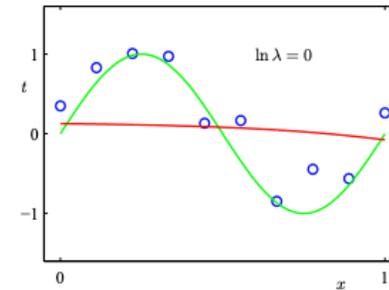
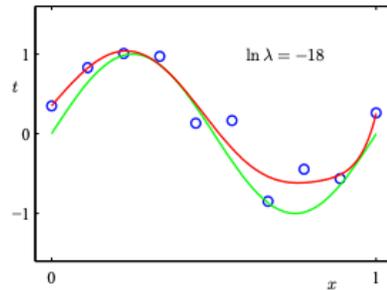
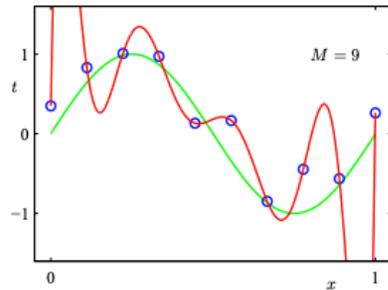
ℓ_2 regularization with non-linear basis: The effect of λ

When we increase regularization coefficient λ :

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
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

ℓ_2 regularization with non-linear basis : A **tradeoff**

when we increase regularization coefficient λ



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?

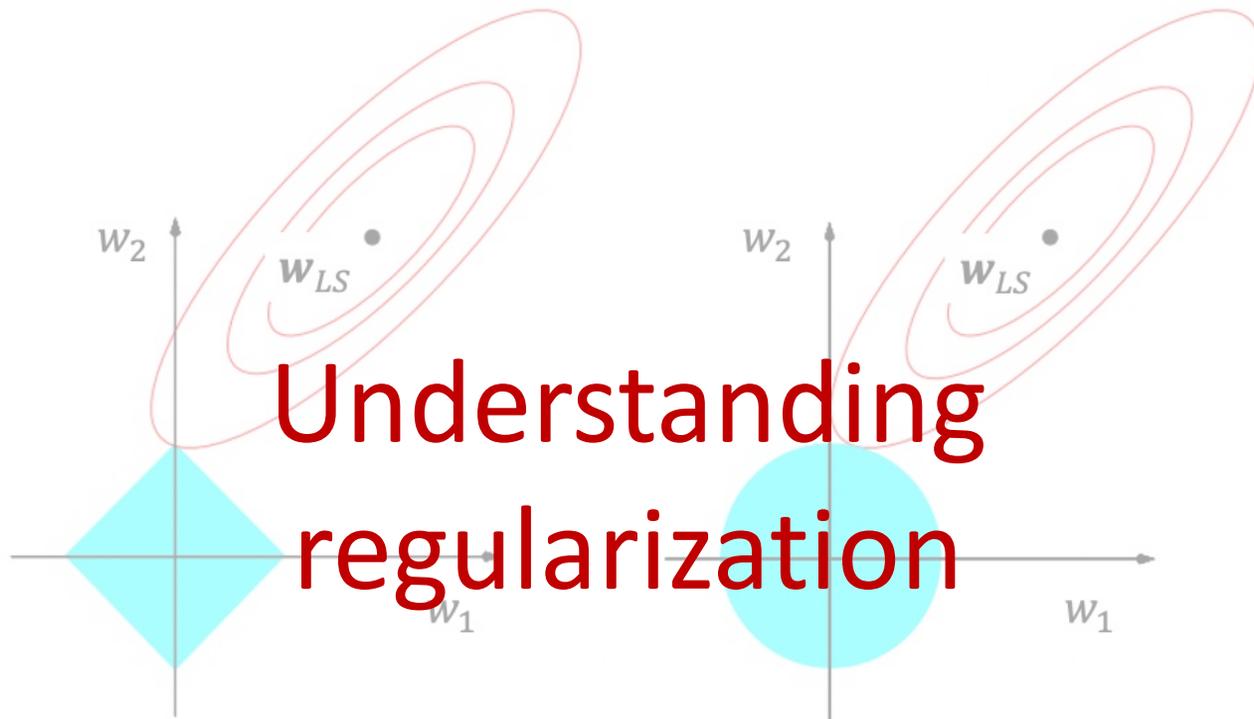
Why is regularization useful?

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).



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

Let's go back to the original linear model.

Simple for ℓ_2 regularization, $\psi(\mathbf{w}) = \|\mathbf{w}\|_2^2: = \mathbf{w}^\top \mathbf{w}$

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

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

Let's go back to the original linear model.

Simple for ℓ_2 regularization, $\psi(\mathbf{w}) = \|\mathbf{w}\|_2^2$:

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

$$\nabla G(\mathbf{w}) = 2(\mathbf{X}^T\mathbf{X}\mathbf{w} - \mathbf{X}^T\mathbf{y}) + 2\lambda\mathbf{w} = 0$$

$$\Rightarrow (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})\mathbf{w} = \mathbf{X}^T\mathbf{y}$$

$$\Rightarrow \mathbf{w}^* = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y}$$

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

Let's go back to the original linear model.

Simple for ℓ_2 regularization, $\psi(\mathbf{w}) = \|\mathbf{w}\|_2^2 = \sum_{i=1}^d w_i^2$

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

$$\nabla G(\mathbf{w}) = 2(\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y}) + 2\lambda \mathbf{w} = 0$$

$$\Rightarrow (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$\Rightarrow \mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

Linear regression with ℓ_2 regularization is also known as **ridge regression**.

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

Regularization when $\mathbf{X}^T \mathbf{X}$ is not invertible: Min-norm solution

When $\mathbf{X}^T \mathbf{X}$ is not invertible $\mathbf{w}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ is not defined.

Regularization when $\mathbf{X}^T \mathbf{X}$ is not invertible: Min-norm solution

When $\mathbf{X}^T \mathbf{X}$ is not invertible $\mathbf{w}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ is not defined.

This could happen when:

1. ∞ many \mathbf{w} s.t. $\mathbf{X}\mathbf{w} = \mathbf{y}$
2. No such \mathbf{w} s.t. $\mathbf{X}\mathbf{w} = \mathbf{y}$

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

Regularization when $\mathbf{X}^T \mathbf{X}$ is not invertible: Min-norm solution

When $\mathbf{X}^T \mathbf{X}$ is not invertible $\mathbf{w}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ is not defined.

This could happen when:

1. ∞ many \mathbf{w} s.t. $\mathbf{X}\mathbf{w} = \mathbf{y}$
2. No such \mathbf{w} s.t. $\mathbf{X}\mathbf{w} = \mathbf{y}$

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

What does L_2 regularization do here?

$$G(\mathbf{w}) = \underbrace{\|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2}_0 + \lambda \|\mathbf{w}\|_2^2$$

consider case 1.

$$\|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = 0 \quad \text{for all } \mathbf{w} \text{ s.t. } \mathbf{X}\mathbf{w} = \mathbf{y}$$

\therefore L_2 regularization chooses \mathbf{w} with smallest $\|\mathbf{w}\|_2^2$ s.t. $\mathbf{X}\mathbf{w} = \mathbf{y}$.

Optional: 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 $\lambda_1 \geq \lambda_2 \geq \cdots \lambda_{D+1} \geq 0$ 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

Optional: 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 ℓ_2 regularization

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

Let's continue with the linear model

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

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

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

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

A “Bayesian view” of ℓ_2 regularization

$$\log(\text{Pr}(y_i | x_i; w)) = \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \frac{(y_i - w^T x_i)^2}{2\sigma^2}$$

(const. (no dependence on w))

Exercise: Assume σ is fixed and given, show that the maximum likelihood estimation for w_* is given by $w_* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{y}$. Here \mathbf{X} is the data matrix with each row corresponding to the features of an example, and \vec{y} is the vector of labels.

Hint: Review the derivation for logistic regression in the last lecture. First write down the probability of seeing the outcomes y_1, \dots, y_n given x_1, \dots, x_n as a function of the value of w_* ; then find the value of w_* that maximizes this probability. You can assume $\mathbf{X}^T \mathbf{X}$ is invertible.

A "Bayesian view" of ℓ_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 the posterior probability (MAP) of w .

Posterior \propto Prior \cdot Likelihood

$$p(w | \text{data}) = \frac{\overset{\text{Prior}}{p(w)} \cdot \overset{\text{Likelihood}}{p(\text{data} | w)}}{p(\text{data})}$$
$$= \frac{p(w \ \& \ \text{data})}{p(\text{data})}$$

A "Bayesian view" of ℓ_2 regularization

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

Bayesian view: A **prior** over w

$$\text{Posterior}(w) = \text{Pr}(w | \text{data}) \propto \prod_{j=1}^d \exp\left(-\frac{w_j^2}{2\gamma^2}\right) \prod_{i=1}^n \exp\left(-\frac{(y_i - w^T x_i)^2}{2\sigma^2}\right)$$

$$\log(\text{Posterior}(w)) = -\frac{\|w\|^2}{2\gamma^2} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w^T x_i)^2 + \text{const.}$$

$\max(\log(\text{Posterior}(w)))$ is same as $\min h(w)$

$$\text{where } h(w) = \|Xw - y\|_2^2 + \frac{\sigma^2}{2\gamma^2} \|w\|^2 \quad (\Psi(w) = \|w\|_2^2)$$

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:

$$\arg \min_{\mathbf{w}} \text{RSS}(\mathbf{w}) \quad \text{subject to } \psi(\mathbf{w}) \leq \beta$$

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

where β is some hyper-parameter.

Finding the solution becomes a *constrained optimization problem*.

Choosing either λ or β can be done by cross-validation.