# COS324:
## INTRODUCTION TO MACHINE LEARNING

**Prof. Yoram Singer**

PRINCETON UNIVERSITY

Topic: Generalization

---

# Thus Far

Definitions of learning problems

Linear and non-linear models

Using differentiable loss for learning

Learning algorithms

---

# Thus Far

Definitions of learning problems

Linear and non-linear models

Using differentiable loss for learning

Learning algorithms

Mentioned in passing through examples **test** loss & error

---

# Thus Far

Definitions of learning problems

Linear and non-linear models
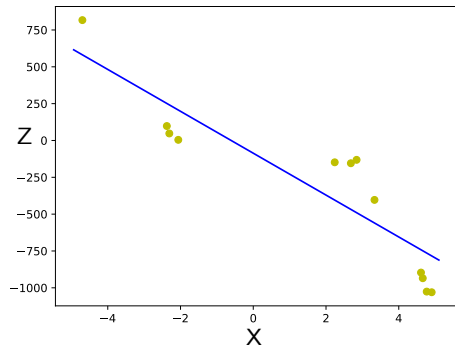
Using differentiable loss for learning

Learning algorithms

Mentioned in passing through examples **test** loss & error

Should the loss/error on unseen data resemble training loss/error ?

## Slide 3

Dataset of examples each has two features $\left\{(x_i, z_i)\right\}_{i=1}^{20}$



Learn a function $f : \mathbf{R} \rightarrow \mathbf{R}$

Regression loss: $\left(f(x) - z\right)^2$

Choose an order $\mathbf{p}$ for a polynomial:

$f(x) = a_0 + a_1 x + a_2 x^2 + \ldots + a_p x^p$

Learn coefficients $a_0, a_1, a_2, \ldots, a_p$

## Learning Polynomials

Replace $x \mapsto \mathbf{x} = (1, x, x^2, x^3, \ldots, x^p)$

For example suppose $x_i = 3$ and $p = 5$ then $x_i \mapsto \mathbf{x}_i = (1, 3, 9, 27, 81, 243)$

## Learning Polynomials

Replace $x \mapsto \mathbf{x} = (1, x, x^2, x^3, \ldots, x^p)$

$$X = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \ldots \\ \mathbf{x}_n \end{bmatrix}$$

For example suppose $x_i = 3$ and $p = 5$ then $x_i \mapsto \mathbf{x}_i = (1, 3, 9, 27, 81, 243)$

## Learning Polynomials

Replace $x \mapsto \mathbf{x} = (1, x, x^2, x^3, \ldots, x^p)$

$$X = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \ldots \\ \mathbf{x}_n \end{bmatrix}$$

For example suppose $x_i = 3$ and $p = 5$ then $x_i \mapsto \mathbf{x}_i = (1, 3, 9, 27, 81, 243)$

| 1 | $x_1$ | $(x_1)^2$ | … | … | $(x_1)^5$ |
|---|-------|-----------|---|---|-----------|
| 1 | $x_2$ | $(x_2)^2$ |   |   | … |
| 1 | $x_3$ | $(x_3)^2$ |   |   | … |
| 1 | $x_4$ | $(x_4)^2$ | … | … | $(x_4)^5$ |

## Learning Polynomials

Replace $x \mapsto \mathbf{x} = (1, x, x^2, x^3, \ldots, x^p)$

$$X = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \ldots \\ \mathbf{x}_n \end{bmatrix}$$

For example suppose $x_i = 3$ and $p = 5$ then $x_i \mapsto \mathbf{x}_i = (1, 3, 9, 27, 81, 243)$

| 1 | $x_1$ | $(x_1)^2$ | ... | ... | $(x_1)^5$ |
|---|---|---|---|---|---|
| 1 | $x_2$ | $(x_2)^2$ | | | ... |
| 1 | $x_3$ | $(x_3)^2$ | | | ... |
| 1 | $x_4$ | $(x_4)^2$ | ... | ... | $(x_4)^5$ |

| $a_1$ |
|---|
| $a_2$ |
| $a_3$ |
| $a_4$ |
| $a_5$ |

---

## Learning Polynomials

Replace $x \mapsto \mathbf{x} = (1, x, x^2, x^3, \ldots, x^p)$

$$X = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \ldots \\ \mathbf{x}_n \end{bmatrix}$$

For example suppose $x_i = 3$ and $p = 5$ then $x_i \mapsto \mathbf{x}_i = (1, 3, 9, 27, 81, 243)$

| 1 | $x_1$ | $(x_1)^2$ | ... | ... | $(x_1)^5$ |
|---|---|---|---|---|---|
| 1 | $x_2$ | $(x_2)^2$ | | | ... |
| 1 | $x_3$ | $(x_3)^2$ | | | ... |
| 1 | $x_4$ | $(x_4)^2$ | ... | ... | $(x_4)^5$ |

| $a_1$ |
|---|
| $a_2$ |
| $a_3$ |
| $a_4$ |
| $a_5$ |

$\approx$

| $z_1$ |
|---|
| $z_2$ |
| $z_3$ |
| $z_4$ |

---

## Learning Polynomials

Replace $x \mapsto \mathbf{x} = (1, x, x^2, x^3, \ldots, x^p)$

$$X = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \ldots \\ \mathbf{x}_n \end{bmatrix}$$

For example suppose $x_i = 3$ and $p = 5$ then $x_i \mapsto \mathbf{x}_i = (1, 3, 9, 27, 81, 243)$

| 1 | $x_1$ | $(x_1)^2$ | ... | ... | $(x_1)^5$ |
|---|---|---|---|---|---|
| 1 | $x_2$ | $(x_2)^2$ | | | ... |
| 1 | $x_3$ | $(x_3)^2$ | | | ... |
| 1 | $x_4$ | $(x_4)^2$ | ... | ... | $(x_4)^5$ |

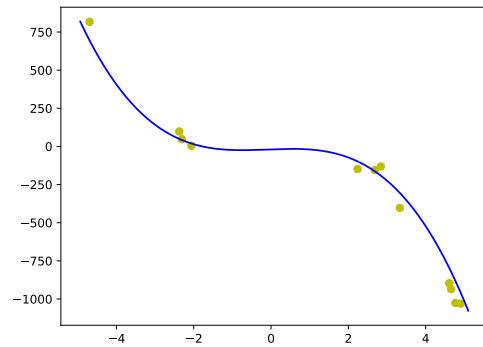| $a_1$ |
|---|
| $a_2$ |
| $a_3$ |
| $a_4$ |
| $a_5$ |

$\approx$

| $z_1$ |
|---|
| $z_2$ |
| $z_3$ |
| $z_4$ |

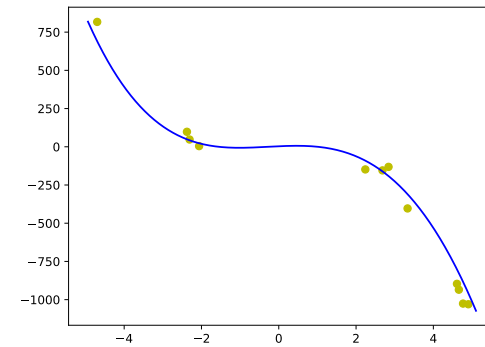$$\min_{\mathbf{a}} \|X\mathbf{a} - \mathbf{z}\|^2$$
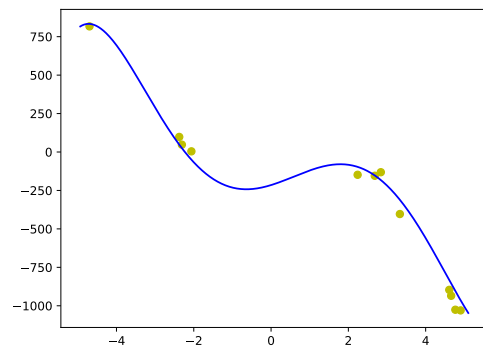
---

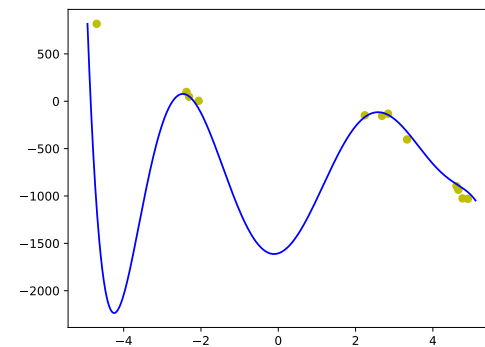## Degree 2 Fit to Training Data

## Degree 3 Fit to Training Data

## Degree 4 Fit to Training Data

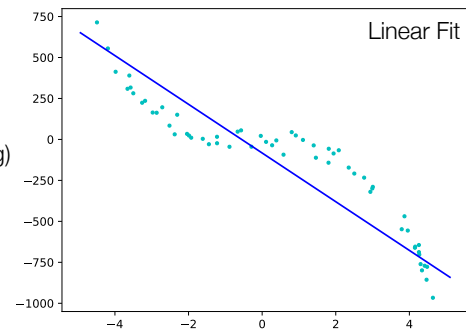## Degree 5 Fit to Training Data
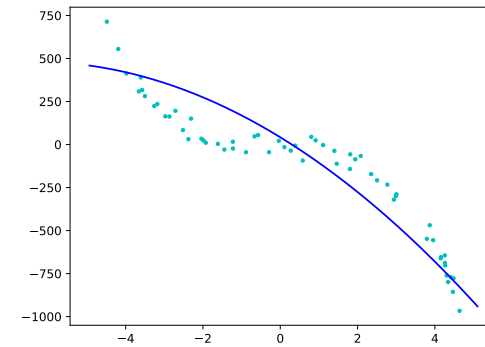
## Degree 7 Fit to Training Data

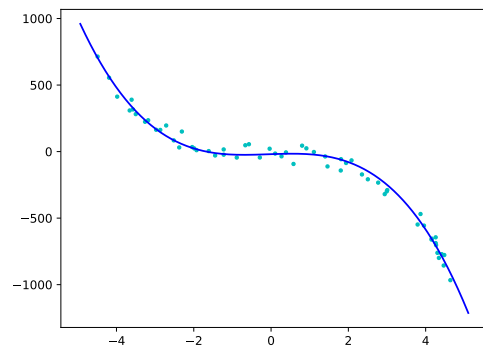## Test Data

Received many more examples

$$\left\{ (x_i, z_i) \right\}_{i=1}^{200}$$

Tested fit on unseen (during training)
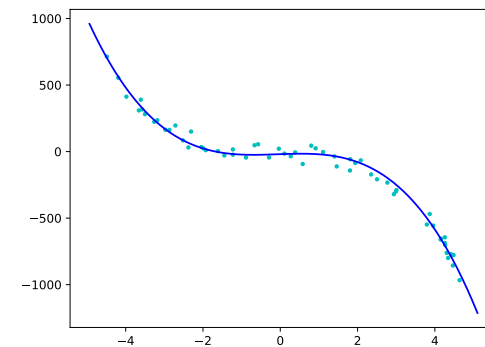
examples



Linear Fit

## Degree 2 Fit to Test Data
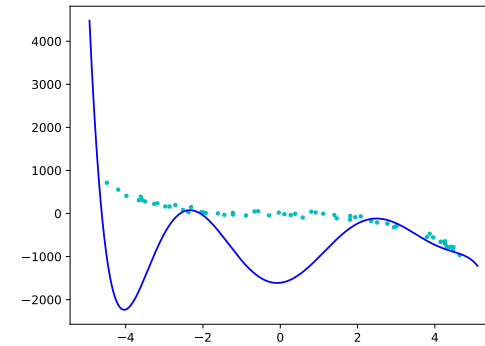
## Degree 3 Fit to Test Data

## Degree 3 Fit to Test Data

Degree 4 Fit to Test Data


Degree 7 Fit to Test Data


What Happened ?


What Happened ?
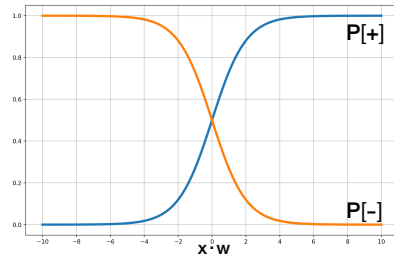
Over Fitting !

## Reminder: Logistic Regression

- Given $\mathbf{x}$ "probability" of y to be +1: $\mathbf{P}\big[+1\,|\,\mathbf{x};\mathbf{w}\big] = \dfrac{1}{1 + e^{-\mathbf{w}\cdot\mathbf{x}}}$

- Probability of y to be -1: $\mathbf{P}\big[-1\,|\,\mathbf{x};\mathbf{w}\big] = 1 - \dfrac{1}{1 + e^{-\mathbf{w}\cdot\mathbf{x}}} = \dfrac{1}{1 + e^{\mathbf{w}\cdot\mathbf{x}}}$

---

## Overfitting in Logistic Regression

Trained 2 logit models:

$$\mathbf{P}\big[y\,|\,\mathbf{x};\mathbf{w_j}\big] = \frac{1}{1 + e^{-y\,\mathbf{w_j}\cdot\mathbf{x}}} \quad j \in [2]$$

Trained with log-loss: for $(\mathbf{x_i}, y_i)$ loss is $-\log\Big(\mathbf{P}\big[y_i\,|\,\mathbf{x_i};\mathbf{w_j}\big]\Big)$
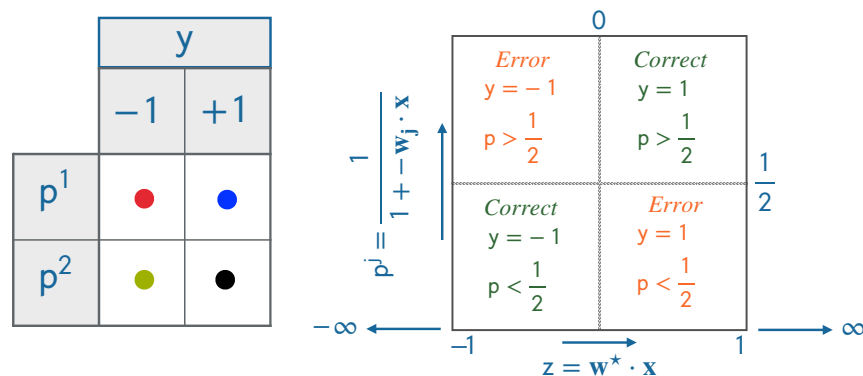
First model was training while guarding for overfitting (more later)

Second model was trained using SGD *without* projections

Predictions: **red** & **blue** first model ; **black** & **yellowish** second model
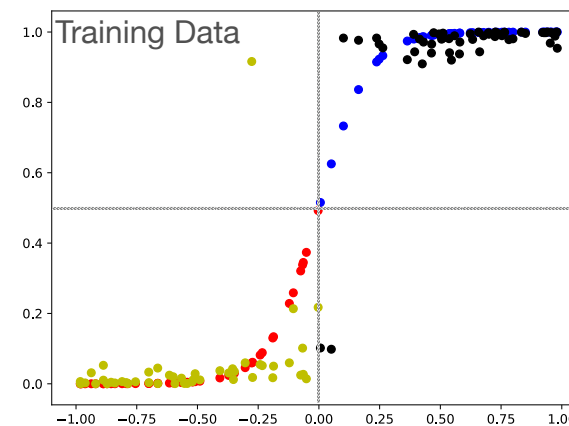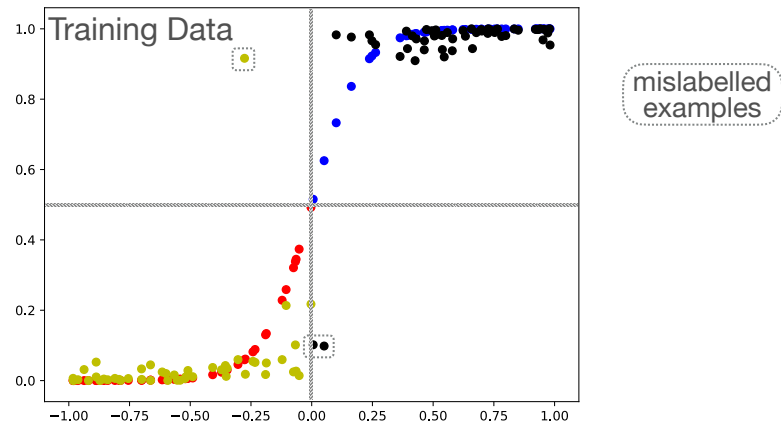
---

## Legend for Graphs

---

## Overfitting in Logistic Regression



Training Data

# Overfitting in Logistic Regression



Training Data

mislabelled examples

# Overfitting in Logistic Regression



Training Data

# Overfitting in Logistic Regression



Test Data

# Overfitting in Classification

Trained a two-layer NN on binary image classification

Thumbnail images: 10x10 (input dimension 100)

Dataset size 10,000

Hidden layer size: 20

Tuned SGD well and ran for many iterations

# Early Stopping

▸ Use a validation set which is not used for training

▸ Check every **k** updates/epochs performance on validation set

▸ Once test-train gap is growing stop training

▸ Works well in practice when scheme is feasible

  ▸ Requires three sets of examples: Train, Validation, Test

  ▸ Loss of stochastic methods not monotone & gap not easy to monitor

Test-Train Gap



Train-Test Gap (large mini-batch)



Train Set

Test Set

Train ∪ Test

## Slide 1

# Underlying Distribution

$$D(\mathbf{x}, \mathbf{y}) =$$
$$D(\mathbf{x})\, D(\mathbf{y} \mid \mathbf{x})$$

## Slide 2

# Underlying Distribution

$$D(\mathbf{x}, \mathbf{y}) =$$
$$D(\mathbf{x})\, D(\mathbf{y} \mid \mathbf{x})$$



$D(\mathbf{x})$

## Slide 3

# Underlying Distribution

$$D(\mathbf{x}, \mathbf{y}) =$$
$$D(\mathbf{x})\, D(\mathbf{y} \mid \mathbf{x})$$



To reason about generalization: Train & test consist of **independent** samples from the **same** distribution

$D(\mathbf{x})$

## Slide 4

# I.I.D Samples

- I.I.D: Identically Independently Distributed

- Generalization analysis typically assumes $\exists D$ :

  unknown distribution $D(\mathbf{x}, y)$

- W.L.O.G assume $\mathbf{x} \in \{0, 1\}^d$   $y \in \{-1, 1\}$

- Identically [no dependence on i]:

  $\forall i \in S : \ D((\mathbf{x}_i, y_i) = (\mathbf{a}, b))$ is $D(\mathbf{a}, b)$

- Independence:

  $D\big((\mathbf{x}_i, y_i) = (\mathbf{a}, b) \ \wedge \ D(\mathbf{x}_i, y_i) = (\mathbf{a}', b')\big) = D(\mathbf{a}, b)\, D(\mathbf{a}', b')$

| $x_0$ | $x_1$ | y | D(x,y) |
|-------|-------|------|--------|
| 0 | 0 | -1 | 0.07 |
| 0 | 0 | 1 | 0.01 |
| 0 | 1 | -1 | 0.03 |
| ... | ... | ... | ... |
| ... | ... | ... | ... |
| 1 | 1 | 1 | 0.005 |

## Generalization Error (deterministic)

Unknown distribution $D(\mathbf{x})$

## Generalization Error (deterministic)

Unknown distribution $D(\mathbf{x})$

Deterministic outcome y given $\mathbf{x}$: $D(y\text{=-}1|\mathbf{x}) = 1$ or $D(y\text{=}1|\mathbf{x}) = 1$

$$\Rightarrow h^{\star}(\mathbf{x}) = \text{sign}(D(y\,|\,\mathbf{x}) - \frac{1}{2})$$

## Generalization Error (deterministic)

Unknown distribution $D(\mathbf{x})$

Deterministic outcome y given $\mathbf{x}$: $D(y\text{=-}1|\mathbf{x}) = 1$ or $D(y\text{=}1|\mathbf{x}) = 1$

$$\Rightarrow h^{\star}(\mathbf{x}) = \text{sign}(D(y\,|\,\mathbf{x}) - \frac{1}{2})$$

Deterministic predictor $f : \{0,1\}^d \to \{-1,1\}$

## Generalization Error (deterministic)

Unknown distribution $D(\mathbf{x})$

Deterministic outcome y given $\mathbf{x}$: $D(y\text{=-}1|\mathbf{x}) = 1$ or $D(y\text{=}1|\mathbf{x}) = 1$

$$\Rightarrow h^{\star}(\mathbf{x}) = \text{sign}(D(y\,|\,\mathbf{x}) - \frac{1}{2})$$

Deterministic predictor $f : \{0,1\}^d \to \{-1,1\}$

Generalization error of $\mathbf{f}$ :

$$\text{err}_D(f) = \sum_{\mathbf{x}} D(\mathbf{x})\,\mathbf{1}[f(\mathbf{x}) \neq h^{\star}(\mathbf{x})]$$

## Generalization Error (deterministic)

Unknown distribution D($\mathbf{x}$)

Deterministic outcome y given $\mathbf{x}$: D(y=-1|$\mathbf{x}$) = 1 or D(y=1|$\mathbf{x}$) = 1

$$\Rightarrow h^\star(\mathbf{x}) = \text{sign}(D(y\,|\,\mathbf{x}) - \frac{1}{2})$$

Deterministic predictor $f : \{0,1\}^d \to \{-1,1\}$

Generalization error of $\mathbf{f}$ :

$$\text{err}_D(f) = \sum_{\mathbf{x}} D(\mathbf{x})\,\mathbf{1}[f(\mathbf{x}) \neq h^\star(\mathbf{x})] = \sum_{\mathbf{x}} D(\mathbf{x}) \sum_{y \in \{-1,1\}} \mathbf{1}[f(\mathbf{x}) \neq y]\,D(y\,|\,\mathbf{x})$$

---

---

$$D(x,y) \qquad D(x,+1) + D(x,-1) = D(x) \qquad D(x,+1) = 1 \;\; or \;\; D(x,-1) = 1$$

$$\sum_{x,y} D(x,y) = 1 \Rightarrow \sum_y \int D(X = x, Y = y)dx$$

$$D(y) = \int D(X = x, Y = y)\, dx$$

$$D(y\,|\,x) = \frac{D(y,x)}{D(x)} = \frac{D(y,x)}{D(x, Y = +\,) + D(x, Y = -\,)}$$

---

# Generalization Error (stochastic)

Unknown distribution $D(\mathbf{x},y)$ & deterministic predictor $f : \{0,1\}^d \rightarrow \{-1,1\}$

# Generalization Error (stochastic)

Unknown distribution $D(\mathbf{x},y)$ & deterministic predictor $f : \{0,1\}^d \rightarrow \{-1,1\}$

Given $\mathbf{x}$ true label $y$ is 1 w.p. $D(y=1|x)$ and -1 w.p. $D(y=-1|x)$

# Generalization Error (stochastic)

Unknown distribution $D(\mathbf{x},y)$ & deterministic predictor $f : \{0,1\}^d \rightarrow \{-1,1\}$

Given $\mathbf{x}$ true label $y$ is 1 w.p. $D(y=1|x)$ and -1 w.p. $D(y=-1|x)$

$f(\mathbf{x})=1 \Rightarrow$ w.p. $D(y=-1|\mathbf{x})$ prediction error ; $f(\mathbf{x})=-1 \Rightarrow$ w.p. $D(y=1|\mathbf{x})$ prediction error

## Generalization Error (stochastic)

Unknown distribution D($\mathbf{x}$,y) & deterministic predictor $f : \{0,1\}^d \to \{-1,1\}$

Given $\mathbf{x}$ true label $y$ is 1 w.p. D(y=1|x) and -1 w.p. D(y=-1|x)

f($\mathbf{x}$)=1 $\Rightarrow$ w.p. D(y=-1|$\mathbf{x}$) prediction error ; f($\mathbf{x}$)=-1 $\Rightarrow$ w.p. D(y=1|$\mathbf{x}$) prediction error

Expected error of $\mathbf{f}$ on $\mathbf{x}$ : $\quad D\big(-f(\mathbf{x})\,|\,\mathbf{x}\big) = \sum_{y\in\{-1,1\}} \mathbf{1}[f(\mathbf{x}) \neq y]\, D(y\,|\,\mathbf{x})$

---

## Generalization Error (stochastic)

Unknown distribution D($\mathbf{x}$,y) & deterministic predictor $f : \{0,1\}^d \to \{-1,1\}$

Given $\mathbf{x}$ true label $y$ is 1 w.p. D(y=1|x) and -1 w.p. D(y=-1|x)

f($\mathbf{x}$)=1 $\Rightarrow$ w.p. D(y=-1|$\mathbf{x}$) prediction error ; f($\mathbf{x}$)=-1 $\Rightarrow$ w.p. D(y=1|$\mathbf{x}$) prediction error

Expected error of $\mathbf{f}$ on $\mathbf{x}$ : $\quad D\big(-f(\mathbf{x})\,|\,\mathbf{x}\big) = \sum_{y\in\{-1,1\}} \mathbf{1}[f(\mathbf{x}) \neq y]\, D(y\,|\,\mathbf{x})$

---

## Generalization Error (stochastic)

Unknown distribution D($\mathbf{x}$,y) & deterministic predictor $f : \{0,1\}^d \to \{-1,1\}$

Given $\mathbf{x}$ true label $y$ is 1 w.p. D(y=1|x) and -1 w.p. D(y=-1|x)

f($\mathbf{x}$)=1 $\Rightarrow$ w.p. D(y=-1|$\mathbf{x}$) prediction error ; f($\mathbf{x}$)=-1 $\Rightarrow$ w.p. D(y=1|$\mathbf{x}$) prediction error

Expected error of $\mathbf{f}$ on $\mathbf{x}$ : $\quad D\big(-f(\mathbf{x})\,|\,\mathbf{x}\big) = \sum_{y\in\{-1,1\}} \mathbf{1}[f(\mathbf{x}) \neq y]\, D(y\,|\,\mathbf{x})$

---

## Generalization Error (stochastic)

Unknown distribution D($\mathbf{x}$,y) & deterministic predictor $f : \{0,1\}^d \to \{-1,1\}$

Given $\mathbf{x}$ true label $y$ is 1 w.p. D(y=1|x) and -1 w.p. D(y=-1|x)

f($\mathbf{x}$)=1 $\Rightarrow$ w.p. D(y=-1|$\mathbf{x}$) prediction error ; f($\mathbf{x}$)=-1 $\Rightarrow$ w.p. D(y=1|$\mathbf{x}$) prediction error

Expected error of $\mathbf{f}$ on $\mathbf{x}$ : $\quad D\big(-f(\mathbf{x})\,|\,\mathbf{x}\big) = \sum_{y\in\{-1,1\}} \mathbf{1}[f(\mathbf{x}) \neq y]\, D(y\,|\,\mathbf{x})$

Generalization error of $\mathbf{f}$ :

$$\text{err}_D(f) = \sum_{\mathbf{x}} D(\mathbf{x}) \sum_{y\in\{-1,1\}} \mathbf{1}[f(\mathbf{x}) \neq y]\, D(y\,|\,\mathbf{x})$$

# Finite Set of Predictors

Suppose we have only **k** predictors — no weight learning: $f_1, \ldots, f_k$

# Finite Set of Predictors

Suppose we have only **k** predictors — no weight learning: $f_1, \ldots, f_k$

One **f** has zero generalization error, rest have generalization error $\geq \epsilon$ :

$$\exists j : \forall (\mathbf{x}, y) : f_j(\mathbf{x}) = h^\star(\mathbf{x}) = y \quad ; \quad \forall i \neq j : \mathrm{err}_\mathbf{D}[f_i(\mathbf{x}) \neq y] \geq \epsilon$$

# Finite Set of Predictors

Suppose we have only **k** predictors — no weight learning: $f_1, \ldots, f_k$

One **f** has zero generalization error, rest have generalization error $\geq \epsilon$ :

$$\exists j : \forall (\mathbf{x}, y) : f_j(\mathbf{x}) = h^\star(\mathbf{x}) = y \quad ; \quad \forall i \neq j : \mathrm{err}_\mathbf{D}[f_i(\mathbf{x}) \neq y] \geq \epsilon$$

Received training set S with only **n** examples sampled independently

# Finite Set of Predictors

Suppose we have only **k** predictors — no weight learning: $f_1, \ldots, f_k$

One **f** has zero generalization error, rest have generalization error $\geq \epsilon$ :

$$\exists j : \forall (\mathbf{x}, y) : f_j(\mathbf{x}) = h^\star(\mathbf{x}) = y \quad ; \quad \forall i \neq j : \mathrm{err}_\mathbf{D}[f_i(\mathbf{x}) \neq y] \geq \epsilon$$

Received training set S with only **n** examples sampled independently

Evaluate errors on S: $\mathrm{err}_S(f_i) = \epsilon_i = \dfrac{1}{n} \sum_{i=1}^{n} \mathbf{1}\big[f_i(\mathbf{x}) \neq y_i\big]$

## Finite Set of Predictors

Suppose we have only $k$ predictors — no weight learning: $f_1, \ldots, f_k$

One $f$ has zero generalization error, rest have generalization error $\geq \epsilon$ :

$$\exists j : \forall (\mathbf{x}, y) : f_j(\mathbf{x}) = h^\star(\mathbf{x}) = y \; ; \quad \forall i \neq j : \mathrm{err}_{\mathbf{D}}[f_i(\mathbf{x}) \neq y] \geq \epsilon$$

Received training set S with only $n$ examples sampled independently

Evaluate errors on S: $\mathrm{err}_S(f_i) = \epsilon_i = \dfrac{1}{n} \sum_{i=1}^{n} \mathbf{1}\big[f_i(\mathbf{x}) \neq y_i\big]$

Choose any $f_j$ for which $\epsilon_j = 0$

---

## Generalization: Finite Case I

Probability that $\epsilon_i = 0$ is at most $(1 - \epsilon)^n \leq e^{-\epsilon n}$ [independence of sample]

---

## Generalization: Finite Case I

Probability that $\epsilon_i = 0$ is at most $(1 - \epsilon)^n \leq e^{-\epsilon n}$ [independence of sample]

Probability $\alpha$ that $\exists i \neq j$ s.t. $\epsilon_i = 0$ is at most $\alpha = (k - 1)\, e^{-\epsilon n}$

---

## Generalization: Finite Case I

Probability that $\epsilon_i = 0$ is at most $(1 - \epsilon)^n \leq e^{-\epsilon n}$ [independence of sample]

Probability $\alpha$ that $\exists i \neq j$ s.t. $\epsilon_i = 0$ is at most $\alpha = (k - 1)\, e^{-\epsilon n}$

If $\alpha \leq \dfrac{1}{k}$ it is unlikely we do not find correct predictor: $(k - 1)\, e^{-\epsilon n} \geq \dfrac{1}{k}$

# Generalization: Finite Case I

Probability that $\epsilon_i = 0$ is at most $(1 - \epsilon)^n \le e^{-\epsilon n}$ [independence of sample]

Probability $\alpha$ that $\exists i \neq j \text{ s.t. } \epsilon_i = 0$ is at most $\alpha = (k-1)\, e^{-\epsilon n}$

If $\alpha \le \dfrac{1}{k}$ it is unlikely we do not find correct predictor: $(k-1)\, e^{-\epsilon n} \ge \dfrac{1}{k}$

This means that we need about $O\!\left(\dfrac{\log(k)}{\epsilon}\right)$ samples

# Generalization: Finite Case II

# Generalization: Finite Case II

Almost always, prefect predictor does not exist: $\epsilon^\star = \min_i \text{err}_D\big(f_i(x)\big) > 0$

# Generalization: Finite Case II

Almost always, prefect predictor does not exist: $\epsilon^\star = \min_i \text{err}_D\big(f_i(x)\big) > 0$

Evaluate errors on S:   $\text{err}_S(f_i) = \epsilon_i = \dfrac{1}{n}\sum_{i=1}^{n} \mathbf{1}\big[f_i(\mathbf{x}) \neq y_i\big]$

## Generalization: Finite Case II

Almost always, prefect predictor does not exist: $\epsilon^\star = \min_i \text{err}_D\big(f_i(x)\big) > 0$

Evaluate errors on S: $\text{err}_S(f_i) = \epsilon_i = \dfrac{1}{n}\sum_{i=1}^{n}\mathbf{1}\big[f_i(\mathbf{x}) \neq y_i\big]$

Choose $f_j$ with the smallest empirical error: $\text{err}_S(f_j)$

---

## Generalization: Finite Case II

Almost always, prefect predictor does not exist: $\epsilon^\star = \min_i \text{err}_D\big(f_i(x)\big) > 0$

Evaluate errors on S: $\text{err}_S(f_i) = \epsilon_i = \dfrac{1}{n}\sum_{i=1}^{n}\mathbf{1}\big[f_i(\mathbf{x}) \neq y_i\big]$

Choose $f_j$ with the smallest empirical error: $\text{err}_S(f_j)$

Generalization error of $f_j$ is: $\text{err}_D\big(f_j(x)\big) = \Delta\epsilon + \epsilon^\star = \Delta\epsilon + \min_i \text{err}_S\big(f_i(x)\big)$

---

## Generalization: Finite Case II

Almost always, prefect predictor does not exist: $\epsilon^\star = \min_i \text{err}_D\big(f_i(x)\big) > 0$
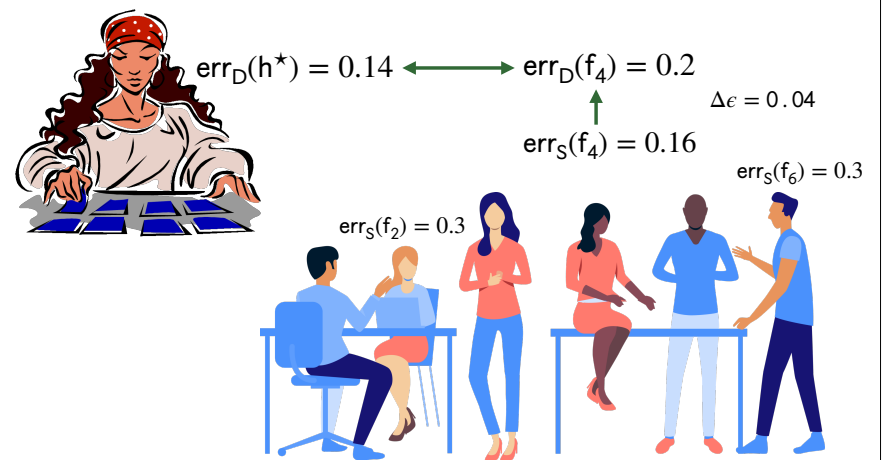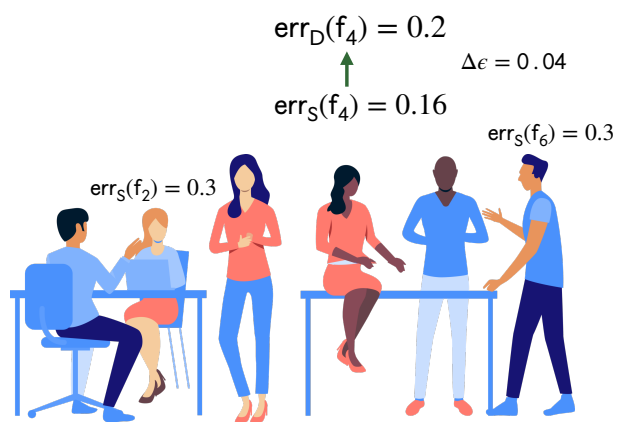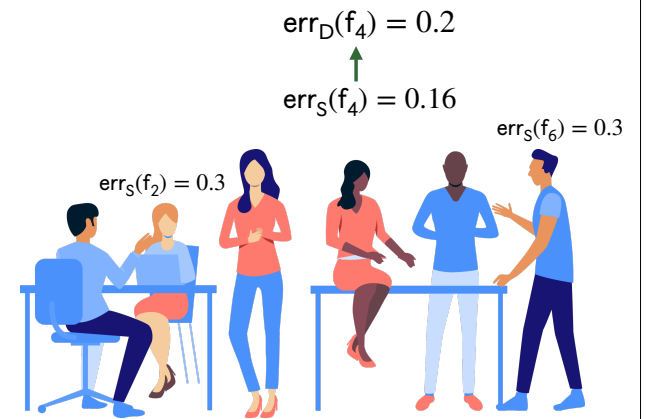
Evaluate errors on S: $\text{err}_S(f_i) = \epsilon_i = \dfrac{1}{n}\sum_{i=1}^{n}\mathbf{1}\big[f_i(\mathbf{x}) \neq y_i\big]$

Choose $f_j$ with the smallest empirical error: $\text{err}_S(f_j)$

Generalization error of $f_j$ is: $\text{err}_D\big(f_j(x)\big) = \Delta\epsilon + \epsilon^\star = \Delta\epsilon + \min_i \text{err}_S\big(f_i(x)\big)$

It takes $O\left(\dfrac{\log(k)}{(\Delta\epsilon)^2}\right)$ sample to get $\Delta\epsilon$-close to f of $\epsilon^\star = \min_i \text{err}_D\big(f_i(x)\big) > 0$

---

## "Continuous Case"

---

## "Continuous Case"

Find best model with weights $\mathbf{w} \in \mathbf{R}^d$

---

## "Continuous Case"

Find best model with weights $\mathbf{w} \in \mathbf{R}^d$

For **bfloat16**: each entry of $\mathbf{w}$ can take $2^{16}$ different values

---

## "Continuous Case"

Find best model with weights $\mathbf{w} \in \mathbf{R}^d$

For **bfloat16**: each entry of $\mathbf{w}$ can take $2^{16}$ different values

Each weight vector corresponds to bit vector of length16d

# "Continuous Case"

Find best model with weights $\mathbf{w} \in \mathbf{R}^d$

For **bfloat16**: each entry of $\mathbf{w}$ can take $2^{16}$ different values

Each weight vector corresponds to bit vector of length16d

We have "only" $2^{16d}$ predictors $\Rightarrow f_1, \ldots, f_{2^{16d}}$

---

# "Continuous Case"

Find best model with weights $\mathbf{w} \in \mathbf{R}^d$

For **bfloat16**: each entry of $\mathbf{w}$ can take $2^{16}$ different values

Each weight vector corresponds to bit vector of length16d

We have "only" $2^{16d}$ predictors $\Rightarrow f_1, \ldots, f_{2^{16d}}$

If $\exists \mathbf{w}^\star$ where $f_{\mathbf{w}^\star}(\mathbf{x}) = y$ for all $\mathbf{x}, y$ with $D(\mathbf{x}, y) > 0$ (0 generalization error):

    it would take only $\tilde{O}(d)$ examples to find it !

---

# "Continuous Case"

Find best model with weights $\mathbf{w} \in \mathbf{R}^d$

For **bfloat16**: each entry of $\mathbf{w}$ can take $2^{16}$ different values

Each weight vector corresponds to bit vector of length16d

We have "only" $2^{16d}$ predictors $\Rightarrow f_1, \ldots, f_{2^{16d}}$

If $\exists \mathbf{w}^\star$ where $f_{\mathbf{w}^\star}(\mathbf{x}) = y$ for all $\mathbf{x}, y$ with $D(\mathbf{x}, y) > 0$ (0 generalization error):

    it would take only $\tilde{O}(d)$ examples to find it !

---

# Caveats

# Caveats

$\tilde{O}(d)$ hides pretty **large** constants

© 2020 Yᴏʀᴀᴍ Sɪɴɢᴇʀ   44

# Caveats

$\tilde{O}(d)$ hides pretty **large** constants

Time of finding $\mathbf{w}^\star$ is **exponential** in $\mathbf{d}$

© 2020 Yᴏʀᴀᴍ Sɪɴɢᴇʀ   44

# Caveats

$\tilde{O}(d)$ hides pretty **large** constants

Time of finding $\mathbf{w}^\star$ is **exponential** in $\mathbf{d}$

Not really learning — exhaustive search

© 2020 Yᴏʀᴀᴍ Sɪɴɢᴇʀ   44

# Caveats

$\tilde{O}(d)$ hides pretty **large** constants

Time of finding $\mathbf{w}^\star$ is **exponential** in $\mathbf{d}$

Not really learning — exhaustive search

Next step:

Incorporate mechanism called regularization into SGD

© 2020 Yᴏʀᴀᴍ Sɪɴɢᴇʀ   44