

Generative Models (Just-so-stories)

- Assume data comes from probabilistic model.
- Reason about outcomes assuming this generative process.

(This is never really true! But it can be instructive)

Example: prediction error.

Assume $y = \vec{w}_*^T \vec{x} + \text{noise}$

Given data, how well can we predict y ?

Let $y_i^{\text{true}} = \vec{w}_*^T \vec{x}_i$

For a learned model \hat{w} , $y_i^{\text{pred}} = \hat{w}^T \vec{x}_i$

prediction error: $\mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n (y_i^{\text{pred}} - y_i^{\text{true}})^2 \right]$

In homework, we showed that for OLS prediction error had expected value $\frac{d}{n}$, provided noise was normally distributed w/ mean zero and variance 1.

OUT OF SAMPLE (generalization)

We don't really care about performance on data we've already seen. What about new data.

$(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)$ what happens on (\vec{x}_{new}, y_{new}) ?
Clearly $\underbrace{\hspace{1cm}}_S$ old and new data have to be related.

Assume (\vec{x}, y) generated by same random process.

PROCESS $\rightarrow (\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n), (\vec{x}_{new}, y_{new})$
I. I. D.

Want to find prediction function

$$f: \mathcal{X} \rightarrow \mathcal{Y} \quad \text{s.t.}$$

$$R[f] \equiv \mathbb{E}_{(x,y)} [\text{loss}(f(x), y)] \quad \text{is small}$$

Can compute

$$R_S[f] = \frac{1}{n} \sum_{i=1}^n \text{loss}(f(x_i), y_i)$$

When is $R[f]$ close to $R_S[f]$?
2

Law of large numbers

If z_1, z_2, \dots, z_n, z i.i.d. then

$$\frac{1}{n} \sum_{i=1}^n z_i \approx \mathbb{E}[z]$$

How close the sample average is to the average is a matter that depends on many properties of z .

Example; biased coin

$$\text{If } z = \begin{cases} 1 & \text{w.p. } p \\ 0 & \text{otherwise (w.p. } 1-p) \end{cases} \quad \begin{matrix} \text{heads} \\ \text{tails} \end{matrix}$$

after n coin flips, expect to see np heads
true number appears to have Gaussian distribution

$$\Pr(\# \text{ heads} \leq np - t) \leq \exp\left(-2 \frac{t^2}{n}\right)$$

$$\Pr(\# \text{ heads} \geq np + t) \leq \exp\left(-2 \frac{t^2}{n}\right)$$

Example: 2 biased coins, $\mathbb{E}[C_1] = p_1$, $\mathbb{E}[C_2] = p_2$

Flip coin 1 and 2 n times, see H_i heads for C_i

H_1, H_2 : Is $p_1 \geq p_2$? How can you tell?

~~XXXXXXXXXXXXXXXXXXXX~~

Assume $p_2 \geq p_1$

$$\begin{aligned} P_r(H_1 \geq H_2) &= P_r(H_1 - H_2 \geq 0) \\ &= P_r(H_1 - H_2 \geq -(p_2 - p_1)n + (p_2 - p_1)n) \\ &\leq \exp\left(-\frac{2(p_2 - p_1)^2 n^2}{2n}\right) = \exp(-n(p_2 - p_1)^2) \end{aligned}$$

If you want to be sure ~~otherwise~~ that the probability of $H_1 \geq H_2$ is less than 0.05 when you assume $p_2 \geq p_1$, then you need

$$n \geq \frac{3}{(p_2 - p_1)^2}$$

closer two alternatives are, the more you need.

Risk $R[f] = \mathbb{E}_{(x,y)} [\text{loss}(f(x), y)]$

Empirical Risk $R_S[f] = \frac{1}{n} \sum_{i=1}^n \text{loss}(f(x_i), y_i)$

For fixed f : $R_S[f] \approx R[f]$

But when is $\min_f R_S[f] \approx \min_f R[f]$?

This is like the biased coin: more possible functions means we need more data.

Idea: Do whatever you'd like on your data S to produce predictor f_S

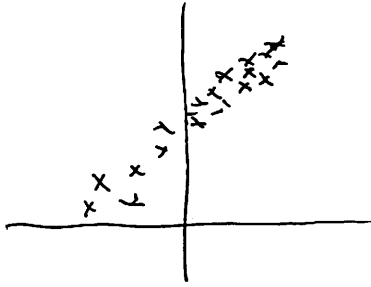
Gather new i.i.d. data S_{NEW} and evaluate

$$R_{S_{\text{NEW}}}[f_S] = \frac{1}{|S_{\text{NEW}}|} \sum_{i \in S_{\text{NEW}}} \text{loss}(f_S(\vec{x}_i), y_i)$$

Example: n example points

\bar{X} features are $\vec{x}_i = \begin{bmatrix} x_i \\ 1 \end{bmatrix}$

\bar{Z} features are $\vec{z}_i = \begin{bmatrix} 1 \\ x_i \\ \vdots \\ x_i^{n-1} \end{bmatrix}$ (square)



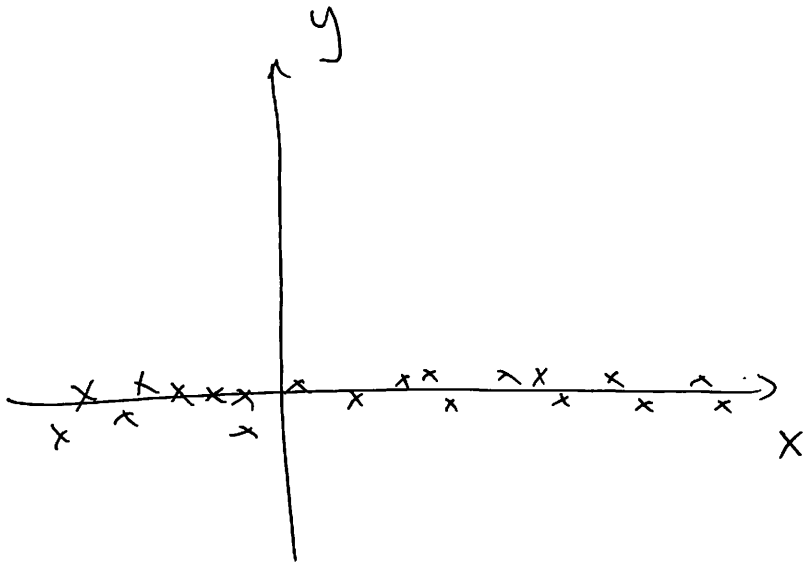
$$\min_{\vec{w}} \|\bar{X} \vec{w} - \vec{y}\| > 0$$

$$\min_{\vec{w}} \|\bar{Z} \vec{w} - \vec{y}\| = 0 !$$

Don't trust training error.

Example: $\arg \min_{\vec{w}} \|\bar{Z} \vec{w} - \vec{y}\|^2 + \gamma \|\vec{w}\|^2$

Is there a value of γ s.t. the fit is best?



$$n = 20$$

Assume $y \sim \mathcal{N}(0, \sigma^2)$
 $\sigma = 0.01$

$$R_S[0] = \frac{1}{n} \sum_{i=1}^n y_i^2$$

Alternative: Let \vec{x} = degree 19 polynomials
 (20 parameters) $\bar{X} = \begin{bmatrix} \vec{x}_1^T \\ \vdots \\ \vec{x}_n^T \end{bmatrix}$

$$\min_{\vec{w}} \|\bar{X} \vec{w} - \vec{y}\|^2 = 0$$

$$\arg \min \rightarrow \vec{w}_{\text{poly}} = (\bar{X}^T \bar{X})^{-1} \bar{X}^T \vec{y}$$

$$R_S[\vec{w}_{\text{poly}}] = 0$$

Simulating new data:

(1) The holdout method. Break data into chunks: (at random!)

TRAIN SET n_T examples

HOLDOUT SET n_H examples

$$n_T + n_H = n$$

Algorithm (Train Set) returns f_T

$$\text{Evaluation (Holdout Set)} = \frac{1}{n_H} \sum_{i \in \text{Holdout}} \text{loss}(f_T(x_i), y_i)$$

(Don't test too many hypotheses on one holdout)

(2) Cross VALIDATION: Break data into F chunks (called folds)

$$S = \{ S_1 \dots S_2 \dots S_F \}$$

for $j=1, \dots, F$:

$$\text{Let train set} = \bigcup_{i \neq j} S_i = S \setminus S_j$$

Algorithm (train set) returns $f_{S \setminus S_j}$

$$\text{Evaluate: } R_j = \frac{1}{|S_j|} \sum_{i \in S_j} \text{loss}(f_{S \setminus S_j}(\vec{x}_i), y_i)$$

$$\text{return } \frac{1}{F} \sum_{j=1}^F R_j$$

intuition, this should have lower variance.

(3) Leave-one-out error (CV to the extreme)

for $k = 1, \dots, n$:

Let train set = $S \setminus \{(\vec{x}_k, y_k)\}$

(delete kth point from training set)

Algorithm (train set) returns $f_{S \setminus k}$

Evaluate $R_k = \text{loss}(f_{S1k}(\vec{x}_k, y_k))$

$$\text{return } \frac{1}{n} \sum_{i=1}^k R_k$$

Example RIDGE REGRESSION

$$\vec{w}_{\text{RIDGE}} = \arg \min_{\vec{w}} \frac{1}{n} \sum_{i=1}^n (\vec{w}^T \vec{x}_i - y_i)^2 + \lambda \|\vec{w}\|^2$$

How do you pick γ ? Choose γ to minimize holdout, cross validation, or leave-one-out error

~~Typically~~ Typically: only look at few γ , ~~logarithmically~~
spaced logarithmically (don't look at all of the γ)

should have $\gamma < \sigma_{\max}(\bar{X})$
 \uparrow data matrix