

GD Gradient Descent

? → Stochastic Gradient Descent

$$\frac{\partial}{\partial \underline{w}} L(D; \underline{w}) = \sum_{i \in D} \frac{\partial}{\partial \underline{w}} l(x_i, y_i; \underline{w})$$

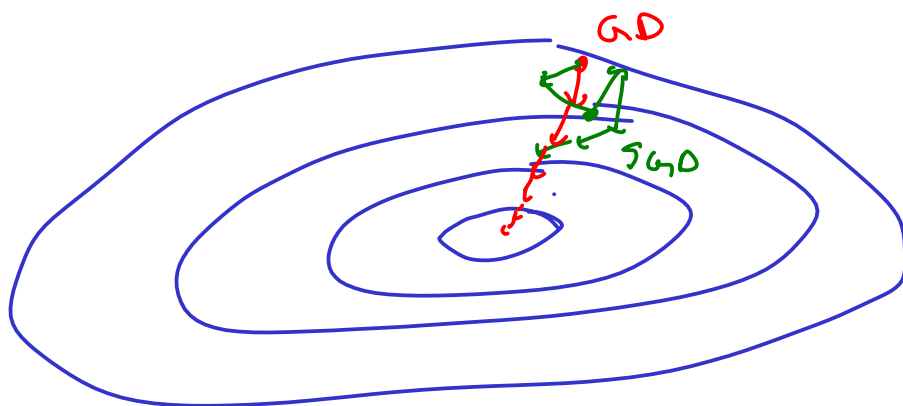
GD is problematic
if $|D|$ is big

while $|\frac{\partial}{\partial \underline{w}} L| > 0.0001$:

$$\underline{w}_{t+1} = \underline{w}_t - \alpha \left[\frac{\partial}{\partial \underline{w}} L(D; \underline{w}) \right]^T$$

Advantages of SGD

- ① Lower memory requirement
- ② Making progress with weights without having to wait full epoch



while
for B in D :
 $\underline{w}_{t+1} = \underline{w}_t - \alpha \sum_{i \in B} \frac{\partial}{\partial \underline{w}} l(x_i, y_i; \underline{w})$
EPOCH
[One iteration over entire Dataset]
Randomly chosen BATCH of data from the dataset D

$|B| = 1$ - SGD

$|B| \approx 64, 32, 128, 256, 512$

Batch SGD

Optimization Perspective on ML

- Data $D = \{(x_i, y_i) \dots\}$
- Select Model $\hat{y}_i = f(x_i; \underline{w})$
- Loss $l(y_i, \hat{y}_i)$

→ Training by GD

$$\underline{w}^* = \arg \min_{\underline{w}} L(D; \underline{w}) = \frac{1}{n} \sum_{i=1}^n l(y_i, \hat{y}_i)$$

→ - New data $x^* \notin D$

Test Data

$$\hat{y}^* = f(x^*; \underline{w}^*)$$

correct?

weight was optimized on the training data

Expected situation data

Training Data \Rightarrow Test Data
(n) R^2 value?

Mathematical

All training and test data samples must be IDENTICALLY distributed.

$$(x_i, y_i) \sim \underset{\downarrow}{\mathbb{P}}(X, Y) \quad \text{Training}$$

$$(x^*, y^*) \sim \mathbb{P}(X, Y) \quad \text{Test / Expected data where your system is supposed to work}$$

In ML, all data must be (i.i.d. assumption)

INDEPENDENT and IDENTICALLY DISTRIBUTED

$$(x_i, y_i) \perp (x_j, y_j)$$

$$(x_i, y_i) \in D, (x_j, y_j) \in D$$

Probabilistic independence $Z_1 \perp Z_2$

$$\textcircled{1} P(Z_1, Z_2) = P(Z_1)P(Z_2) \Leftrightarrow \textcircled{2} P(Z_1 | Z_2) = P(Z_1)$$

Why do we need Independence assumption?

Optimization $\xleftrightarrow{\text{Independence assumption}}$ Probabilistic perspective

$$\arg \max_{\underline{w}} P(D | \underline{w})$$

maximum likelihood estimator

$$P(D | \underline{w})$$

$$= P((x_1, y_1) \dots (x_n, y_n) | \underline{w})$$

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

$$\arg \max_{\underline{w}} \prod_{i=1}^n P(x_i, y_i | \underline{w})$$

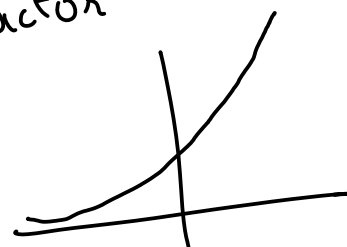
$$P(z_1, z_2)$$

$$= P(z_1)P(z_2)$$

$$\arg \min_{\underline{w}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{y}_i; \underline{w})$$

$$P(x_i, y_i | \underline{w}) = \frac{1}{Z} \exp(-\ell(y_i, \hat{y}_i; \underline{w}))$$

normalizing factor



$$= \arg \max_{\underline{w}} - \frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{y}_i; \underline{w})$$

exponential is a monotonically increasing function

$$= \arg \max_{\underline{w}} \exp\left(-\frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{y}_i; \underline{w})\right)$$

$$= \arg \max_{\underline{w}} \prod_{i=1}^n \exp\left(-\frac{1}{n} \ell(y_i, \hat{y}_i; \underline{w})\right)$$

→ i.i.d assumption

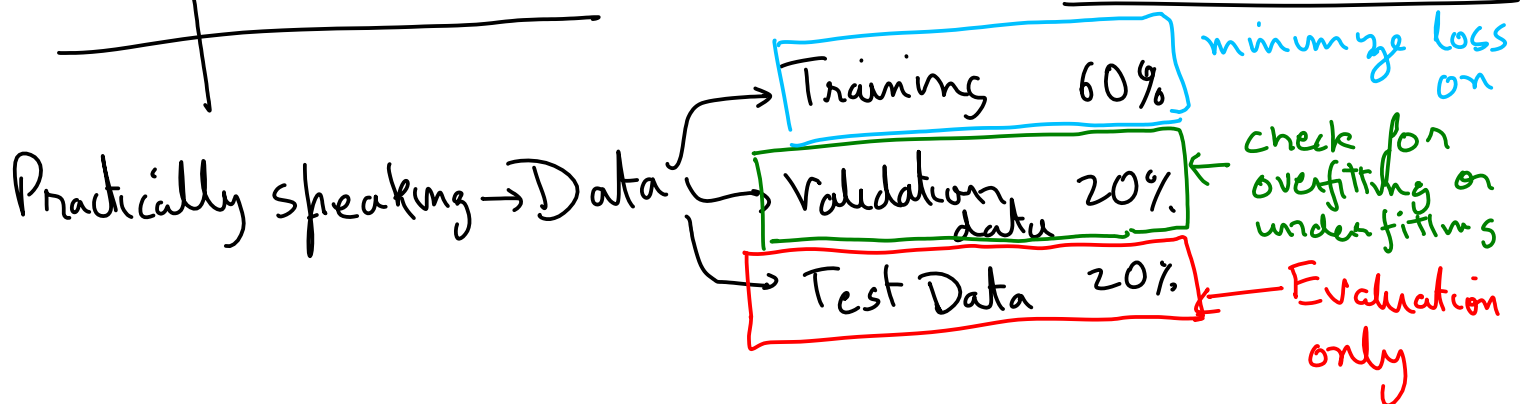
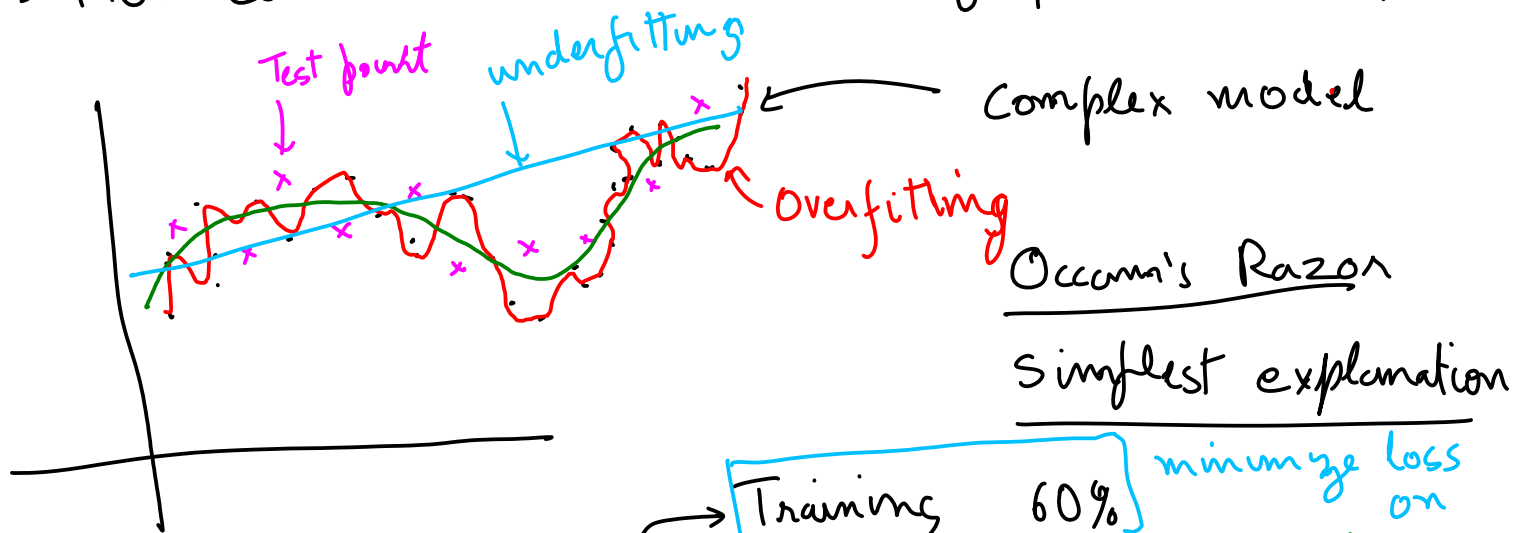
Beyond Linear Models:

→ Does minimizing the training loss
minimize the test loss?

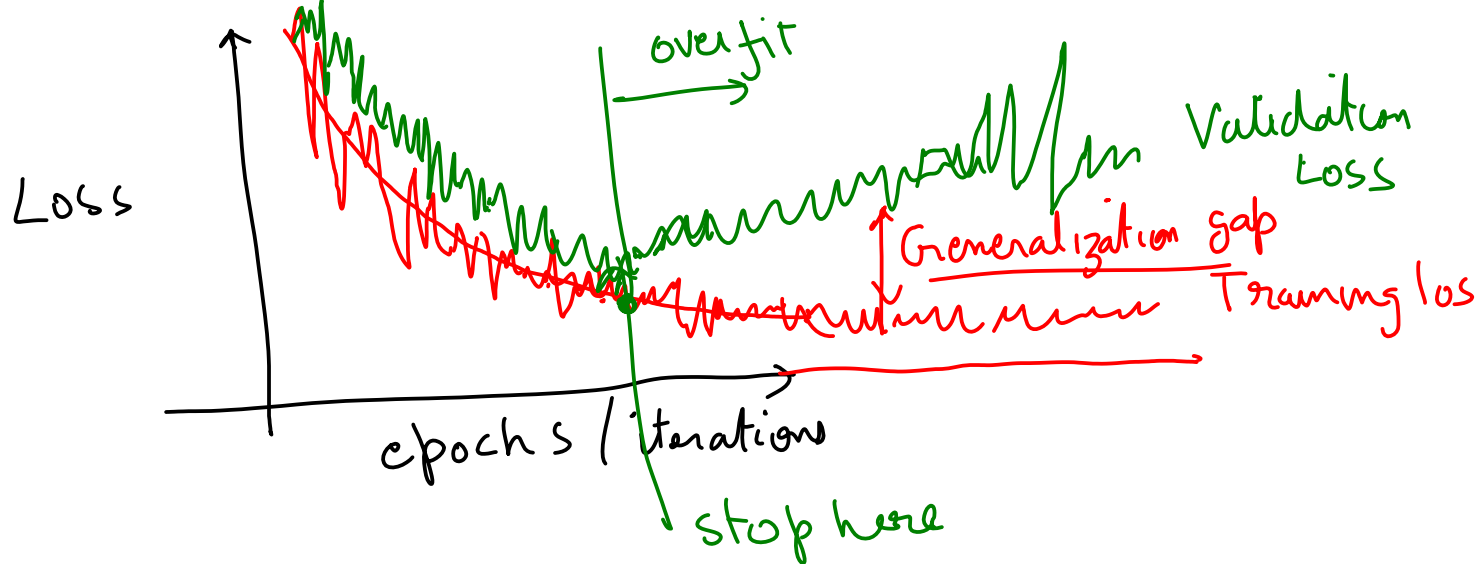
a) Always

b) Sometimes → When it does / When it does not?
(Overfitting)

→ How can we ensure that the gap is small?



Overfitting means Test loss \gg Training loss
detected using Validation loss \gg Training loss



Early stopping technique

Expectation $\rightarrow E[X] = \int_{-\infty}^{\infty} x f(X=x) dx = \sum_x x P(X=x)$

Sample mean $\rightarrow \frac{1}{n} \sum_{i=1}^n X_i$