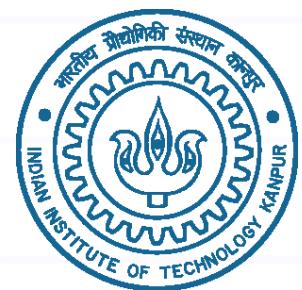


# Neural Networks Optimization

---

*EE698V - Machine Learning for Signal Processing*

Vipul Arora



# Optimization Schemes

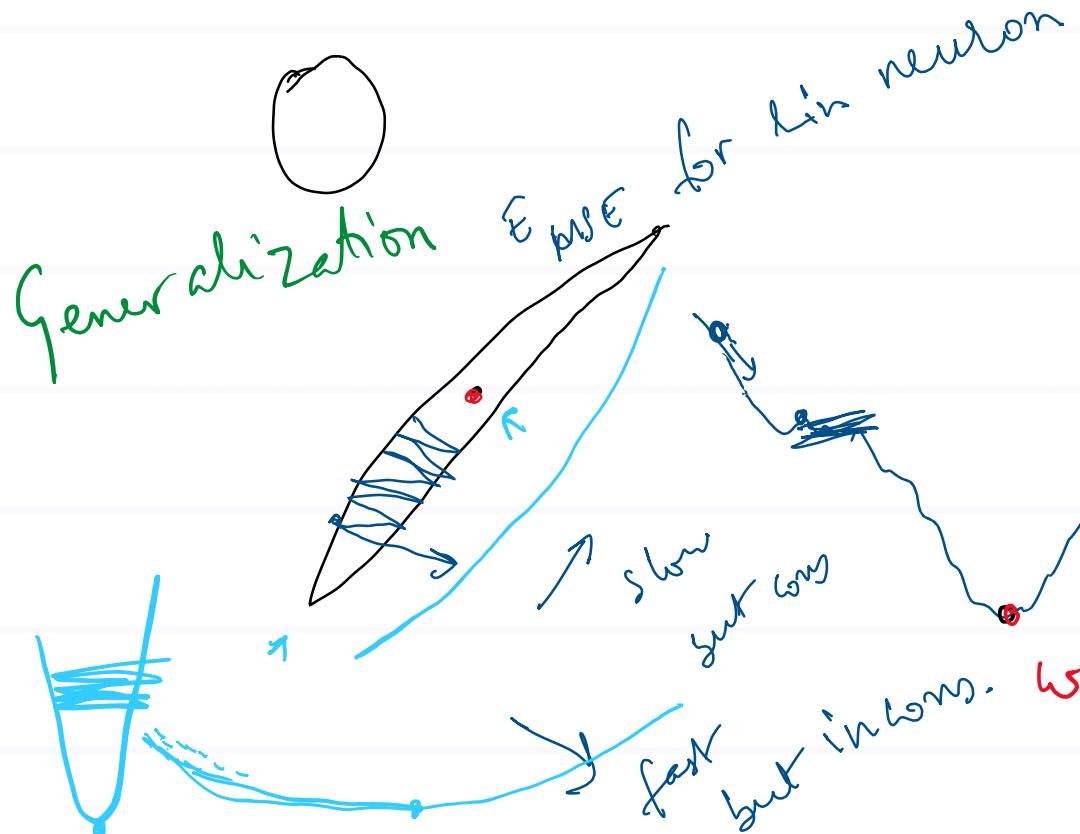
References:

- Hinton's slides: [https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture\\_slides\\_lec6.pdf](https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf)
- Hinton's video lectures: [6a](#), [6b](#), [6c](#), [6d](#), [6e](#)

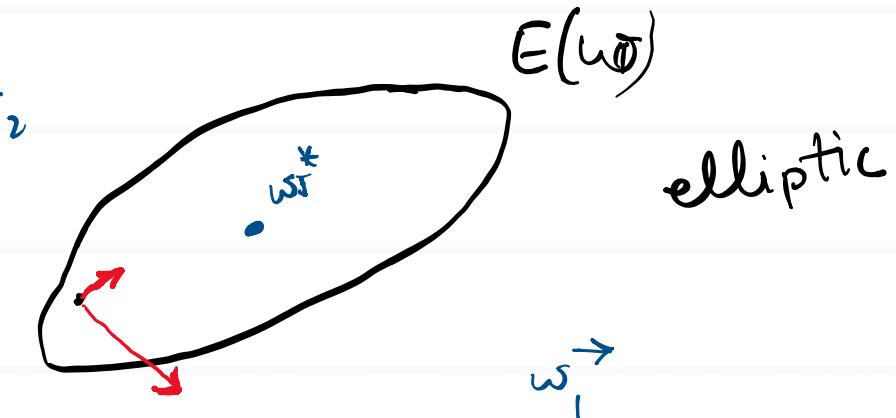
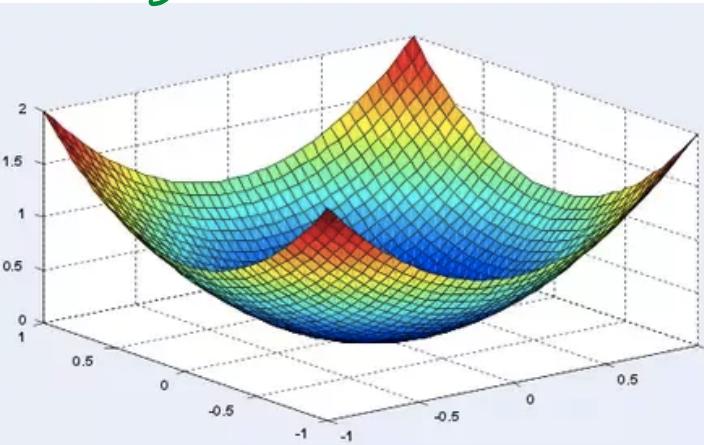
# Optimization Schemes

$$\begin{aligned} E &= \|t - y(x, w)\|^2 \\ &= \mathcal{L}(t, y(x, w)) \end{aligned}$$

- Gradient Descent



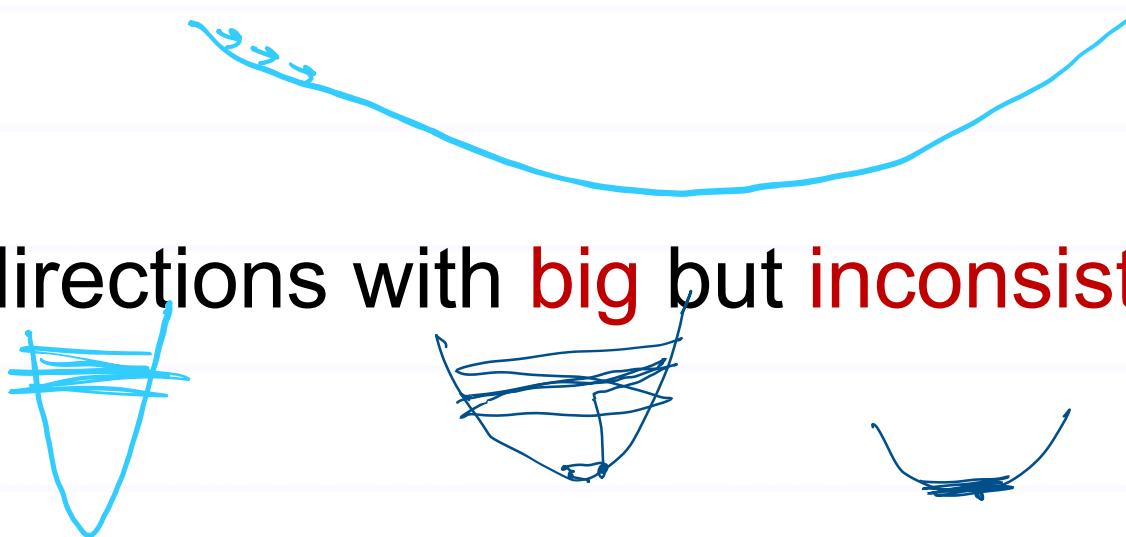
$$w \leftarrow w - \eta \frac{\partial E}{\partial w}$$



# Learning Rate

---

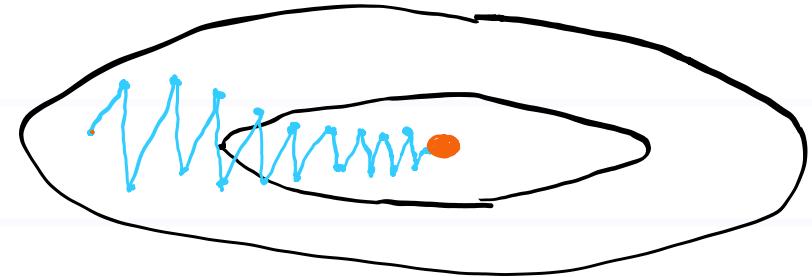
- Move **quickly** in directions with **small** but **consistent** gradients.
- Move **slowly** in directions with **big** but **inconsistent** gradients.
- Towards the **end** of learning, **reduce** the learning rate



# Momentum in Gradient Descent

$$v \leftarrow \beta v + (-\beta) \Delta w$$

vel      in.vel.  
 $v = u + at$       acc.



$$w = w - \alpha v$$

$\Delta w$

$$\sim (s = ut + \frac{1}{2}at^2)$$

# RmsProp

---

- DIY
- References:
  - <https://towardsdatascience.com/a-look-at-gradient-descent-and-rmsprop-optimizers-f77d483ef08b>
  - Hinton's slides:  
[https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture\\_slides\\_lec6.pdf](https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf)
  - Hinton's video lectures: [6a](#), [6b](#), [6c](#), [6d](#), [6e](#)

leaky ReLU

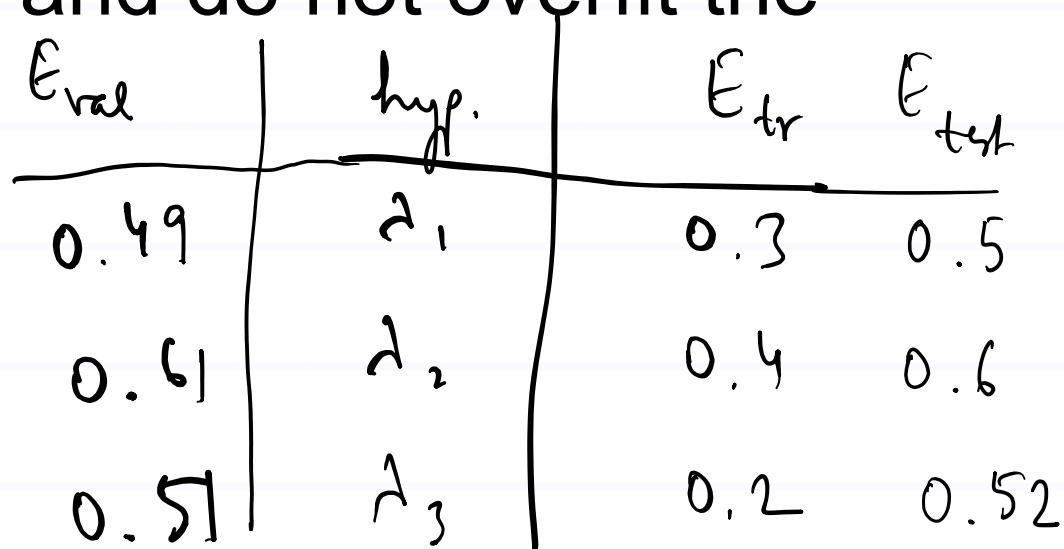
$$\alpha < 0 \quad y = \alpha x \quad y = x$$

# Hyperparameter Tuning

- model arch.
  - layers
  - neurons / l
  - $\eta$
  - $\lambda$

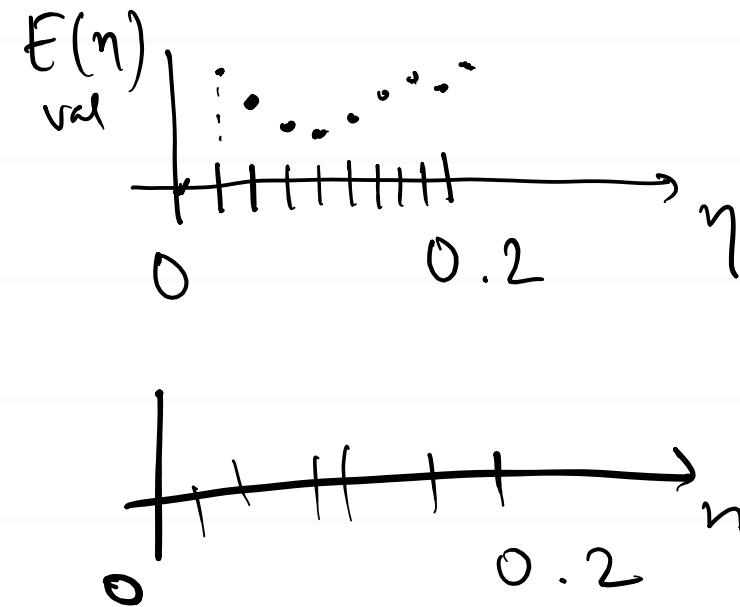
any hyper param.

- How to find the best learning rate?
- Do not optimize it on the test set and do not overfit the training set
- Thus, use validation set
- **Reference:**
  - DeepLearningBook – Section 7.8

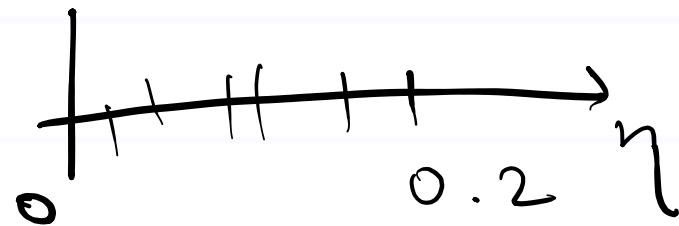


# Hyperparameter Search

- Grid search vs random search



grid search



$\eta \sim \text{unif}(0, 0.2)$

# Validation

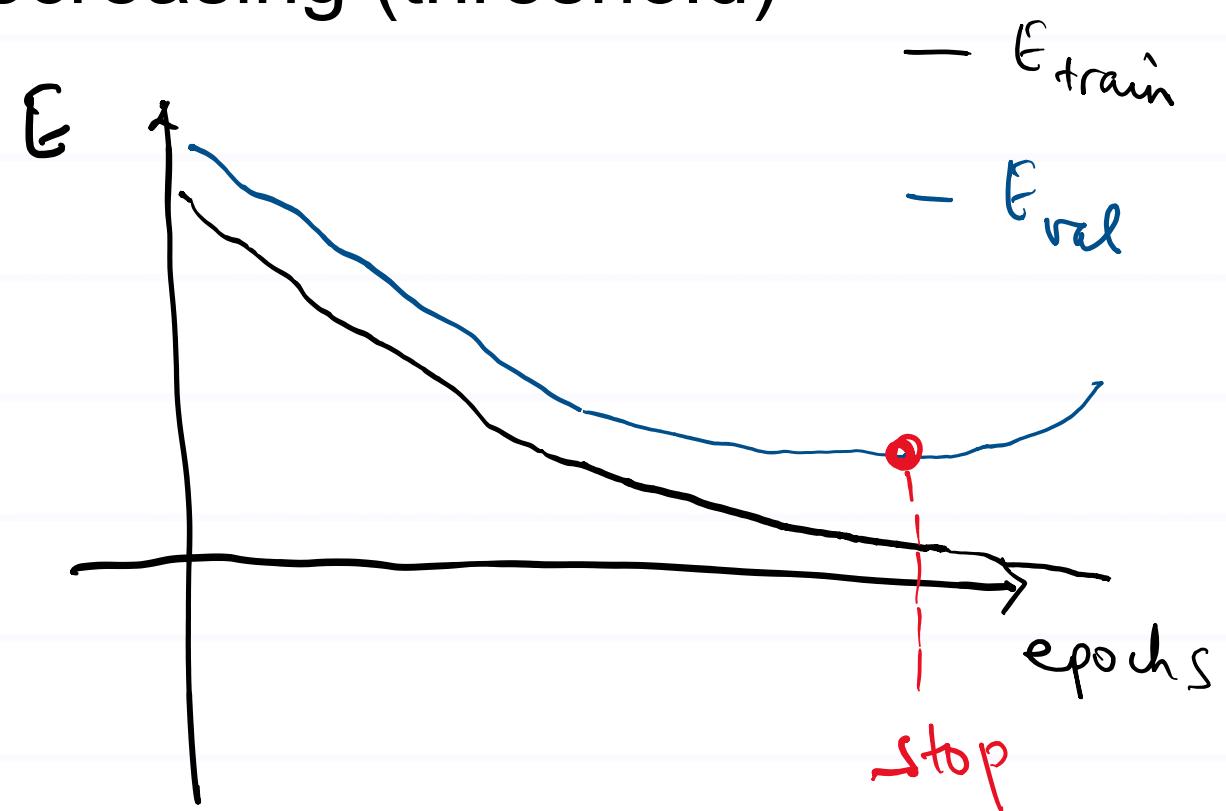
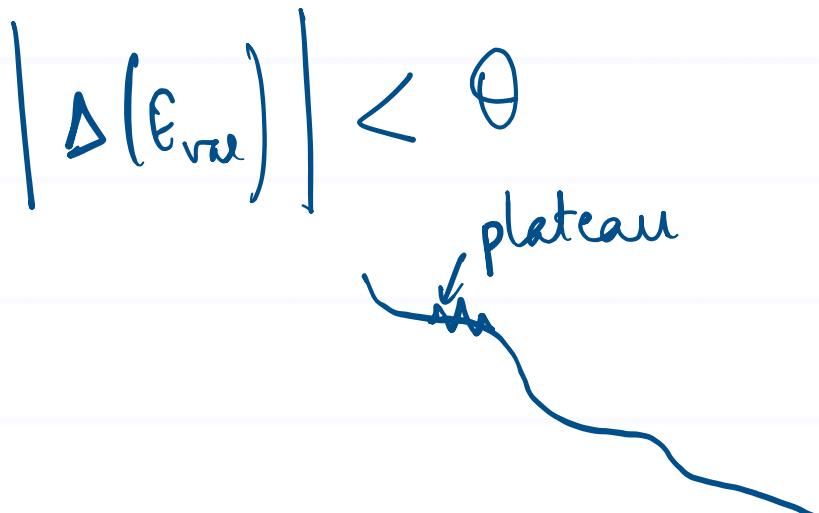
---

- For efficient implementation
  - use separate machine
  - reduce data size
  - reduce frequency of checks

Epoch → seeing all training samples once

# Stopping criterion

- When training error stops decreasing (threshold)
- Monitor N epochs



# Weight Initialization

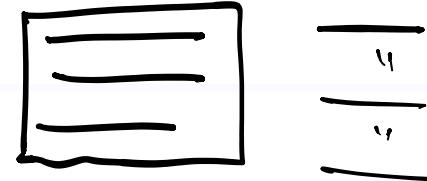


- Random weights

$$h_{i,j} = \sum_{l=0} w_{i,j,l} x_l$$

$(0,1)$

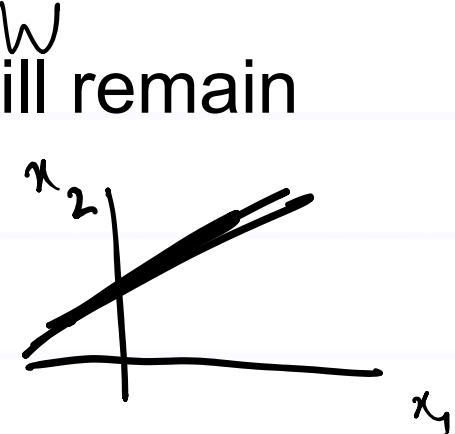
$$w_{i,j,l} = w_{i,j,l}$$



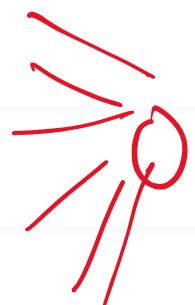
- Same weights will update in the same way and will remain same

- Small weights

$$w \sim 2.1, 2.2, 2.1, 4.3, \dots$$



- more fan-in => overshoot in training



# Data Normalization

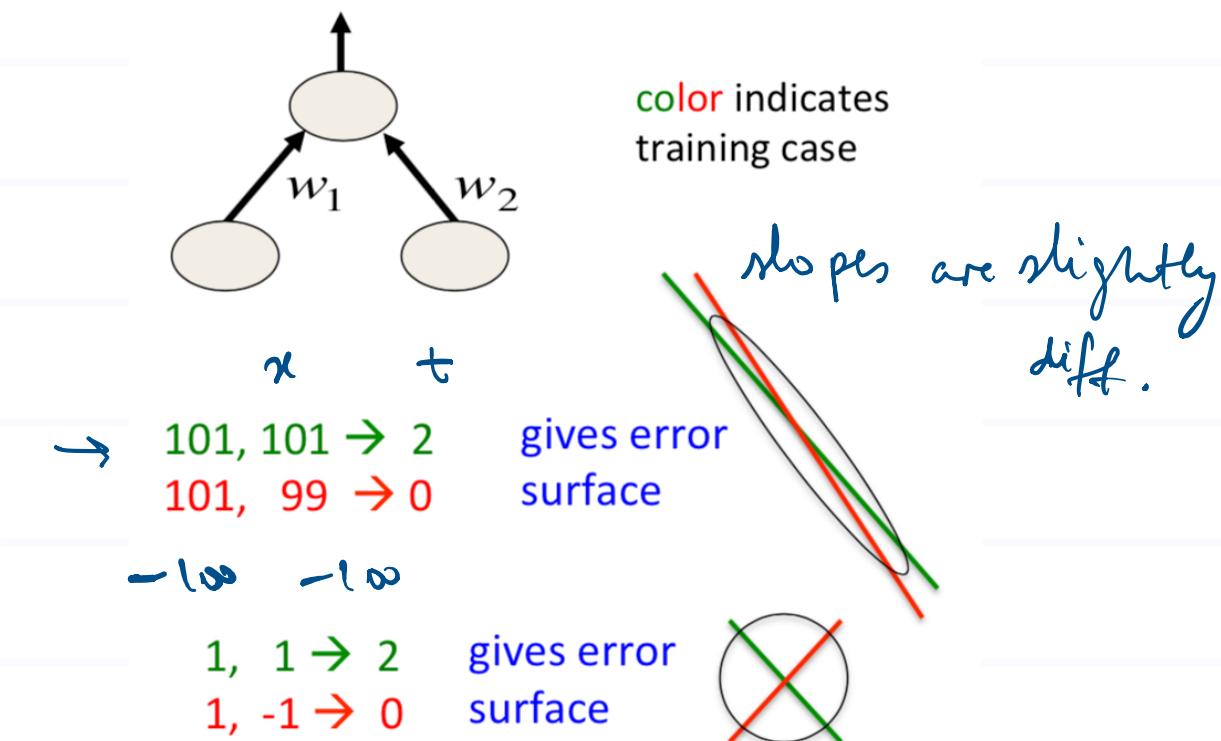
- Shifting of data helps

- make it zero mean

- $x_{i,s} \leftarrow x_{i,s} - \mu_i$

$$w_1|01| + w_2|01| = 2$$

$$w_1|01| + w_2|99| = 0$$



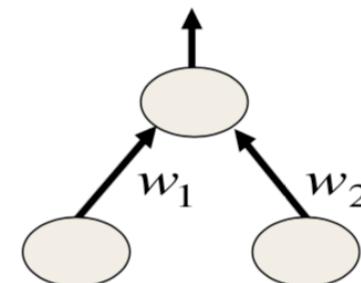
# Data Normalization

$$\frac{w_1 x + \dots + w_{19} x^{19}}{127} \rightarrow \frac{0,127}{127^{19}}$$

- Scaling of data helps

- different variance of different dimensions
- make each dim unit variance
- $x_{i,s} \leftarrow x_{i,s}/\sigma_i$

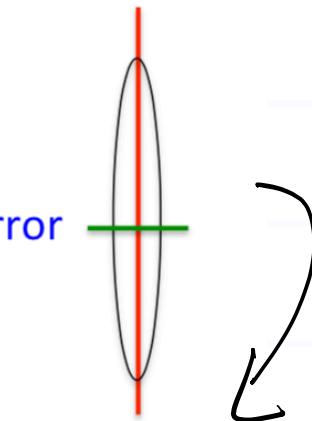
Overall,  $x_{i,s} \leftarrow (x_{i,s} - \mu_i)/\sigma_i$



color indicates weight axis

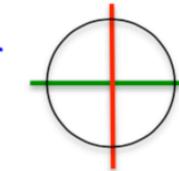
$$0.1, \textcolor{red}{10} \rightarrow 2 \\ 0.1, \textcolor{green}{-10} \rightarrow 0$$

gives error surface



$$1, \textcolor{red}{1} \rightarrow 2 \\ 1, \textcolor{green}{-1} \rightarrow 0$$

gives error surface



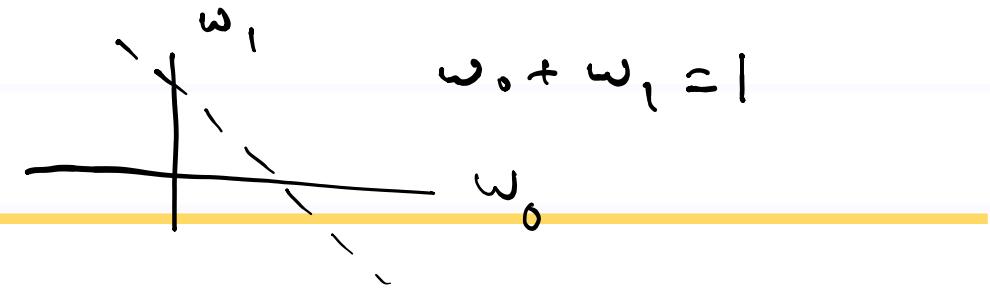
# Regularization

References:

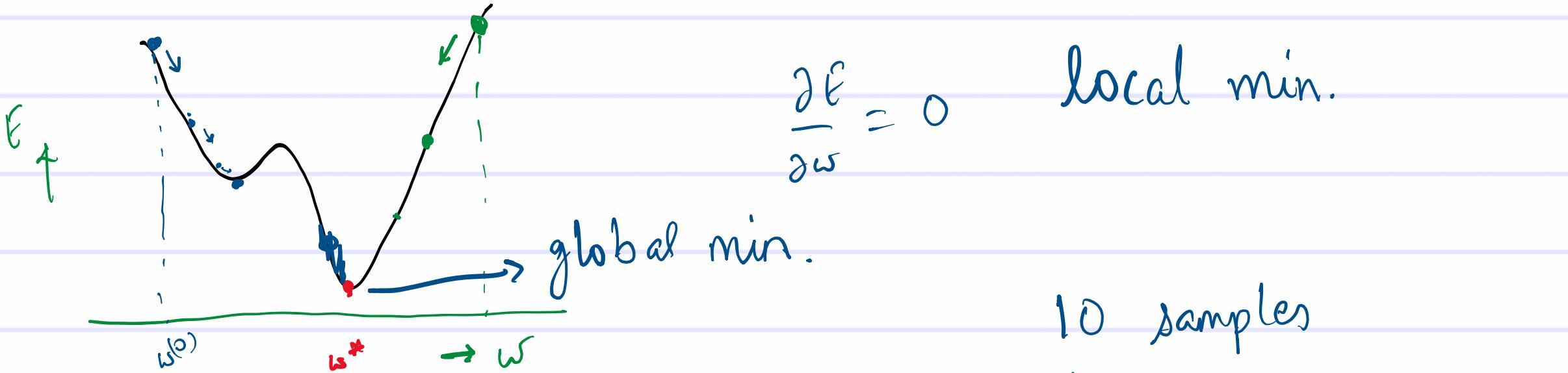
- Deep Learning Book, Chapter 7

# Regularization

---



- Constrain weights
- Prior knowledge injection
- Prefer simpler model classes
- Make underdetermined problems determined
- Ensemble methods



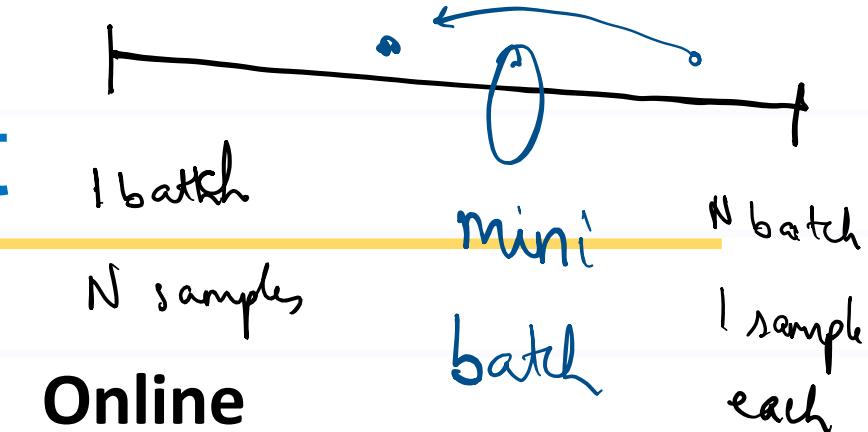
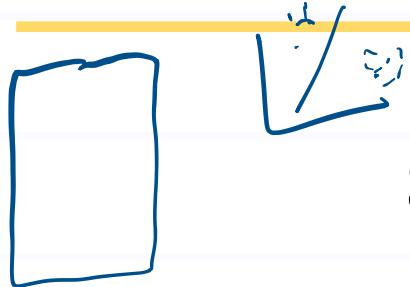
$$\frac{\partial E}{\partial w} = 0$$

local min.

global min.

- $w^{(0)}$
  - 10 samples - find  $\Delta w$
  - $w^{(1)} \leftarrow w^{(0)} + \Delta w$
- $w^{(0)}$
- 5 samples randomly, find  $\Delta w$
- $w^{(0,1)} \leftarrow w^{(0)} + \Delta w$
- 5 samples
- 10 samples
- $w^{(1)} \leftarrow w^{(0,1)} + \Delta w$

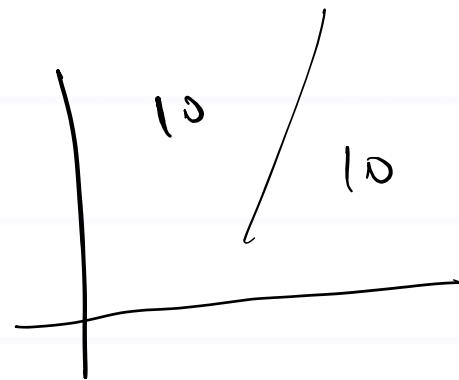
# Stochastic Gradient Descent



- Computationally efficient (matrix, GPU) and takes less time for one epoch
- Redundancy: less improvement in one epoch
- Stable at the minima
- Computationally inefficient and time consuming for one epoch
- Each sample improves weights
- Gets too noisy near the minima

Thus, use randomly sampled **mini-batches**

*One strategy:* Start with small batch size and increase it with epochs



Online

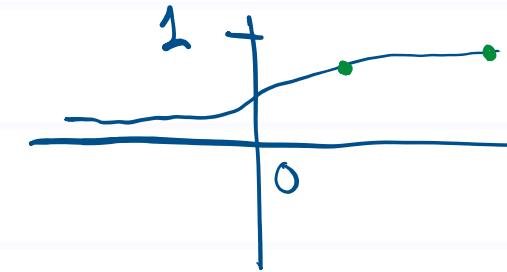
# Adding Noise improves Generalization

---

- Output Noise

- $\{0,1\} \rightarrow \left\{ \frac{\epsilon}{1-K}, 1 - \epsilon \right\}$

Label smoothing



# Adding Noise improves Generalization

---

- Noise in weights
  - Add a gaussian noise to weights
  - $\mathbf{w} \leftarrow \mathbf{w} + \epsilon_{\mathbf{w}}$  where  $\epsilon_{\mathbf{w}} \sim \mathcal{N}(0, \sigma^2 I)$
  - This is equivalent to adding a regularization term to the loss function

# Early Stopping

---

- Limits the search region volume to  $\propto \eta\tau$

↑  
Step  
size  
# of epochs

# Ensemble models or Bagging

---

- Final output is the average of many models
- Not all models make same errors on the test set

# Dropout

forward pass

- Mask neurons randomly

for 1 batch



# Other References

---

- <https://www.deeplearningbook.org/contents/optimization.html>

# Keras library for Python

---

- <https://keras.io/>
- <https://keras.io/#getting-started-30-seconds-to-keras>

# Announcements

---

- Assignment 3 deadline extended by 1 day
- Quiz 2 next week (Fri or Sat?)
- Quiz 1 marks today
  - Copying cases caught
    - will deduct 5% marks from total
    - more severe penalty next time
- Mid sem exam will be written as well as coding based

