# Lecture 5 recap

Prof. Leal-Taixé and Prof. Niessner

# Neural Network



input layer

hidden layer 1    hidden layer 2    hidden layer 3

output layer

Width

Depth

# Gradient Descent for Neural Networks



$h_0$

$h_1$

$x_0$

$x_1$

$x_2$

$h_2$

$h_3$

$y_0$ ← $t_0$

$y_1$ ← $t_1$

input layer

hidden layer

output layer

$$\nabla_{w,b} f_{\{x,t\}}(w) = \begin{bmatrix} \dfrac{\partial f}{\partial w_{0,0,0}} \\ \cdots \\ \cdots \\ \dfrac{\partial f}{\partial w_{l,m,n}} \\ \cdots \\ \cdots \\ \dfrac{\partial f}{\partial b_{l,m}} \end{bmatrix}$$

$$L_i = (y_i - t_i)^2$$

$$h_j = A(b_{0,j} + \sum_k x_k w_{0,j,k})$$

$$y_i = A(b_{1,i} + \sum_j h_j w_{1,i,j})$$

Just simple: $A(x) = \max(0, x)$

# Stochastic Gradient Descent (SGD)

$$\theta^{k+1} = \theta^k - \alpha \nabla_\theta L(\theta^k, x_{\{1..m\}}, y_{\{1..m\}})$$

$$\nabla_\theta L = \frac{1}{m} \sum_{i=1}^{m} \nabla_\theta L_i$$

$k$ now refers to $k$-th iteration

$m$ training samples in the current batch

Gradient for the $k$-th batch

Note the terminology: iteration vs epoch

# Gradient Descent with Momentum

$$v^{k+1} = \beta \cdot v^k + \nabla_\theta L(\theta^k)$$

accumulation rate
('friction', momentum)

velocity

Gradient of current minibatch
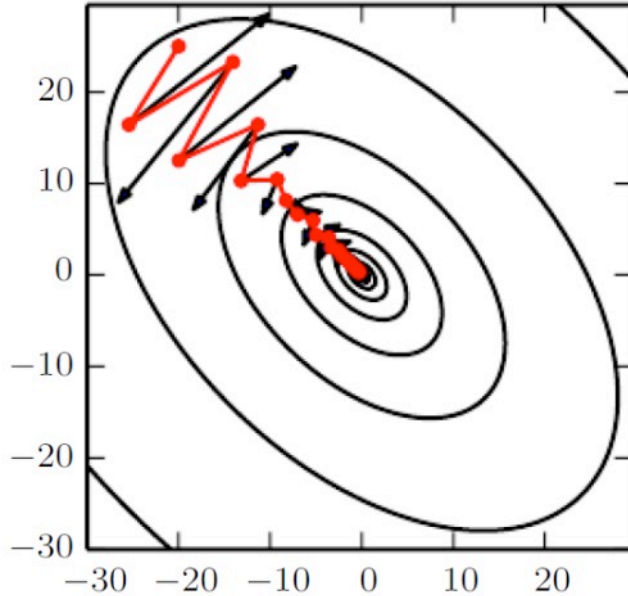
$$\theta^{k+1} = \theta^k - \alpha \cdot v^{k+1}$$

model

learning rate

velocity

Exponentially-weighted average of gradient

Important: velocity $v^k$ is vector-valued!

# Gradient Descent with Momentum



Step will be largest when a sequence of gradients all point to the same direction

Hyperparameters are $\alpha, \beta$
$\beta$ is often set to 0.9

$$\theta^{k+1} = \theta^k - \alpha \cdot v^{k+1}$$

Fig. credit: I. Goodfellow

# RMSProp

$$s^{k+1} = \beta \cdot s^k + (1 - \beta)[\nabla_\theta L \circ \nabla_\theta L]$$

Element-wise multiplication

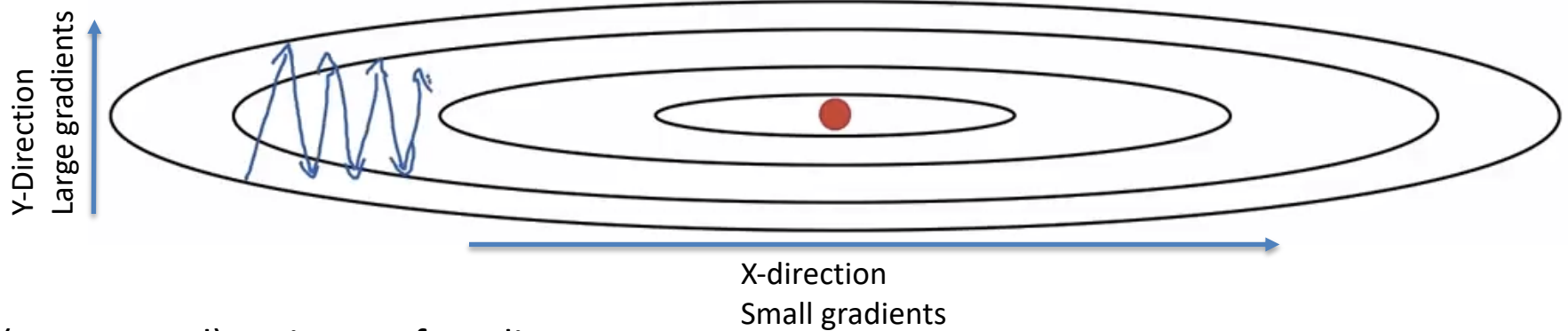$$\theta^{k+1} = \theta^k - \alpha \cdot \frac{\nabla_\theta L}{\sqrt{s^{k+1}} + \epsilon}$$

Hyperparameters: $\alpha$, $\beta$, $\epsilon$

Needs tuning!

Often 0.9

Typically $10^{-8}$

# RMSProp



Y-Direction
Large gradients

X-direction
Small gradients

(uncentered) variance of gradients
-> second momentum

$$s^{k+1} = \beta \cdot s^k + (1 - \beta)[\nabla_\theta L \circ \nabla_\theta L]$$

$$\theta^{k+1} = \theta^k - \alpha \cdot \frac{\nabla_\theta L}{\sqrt{s^{k+1}} + \epsilon}$$

We're dividing by square gradients:
- Division in Y-Direction will be large
- Division in X-Direction will be small

Can increase learning rate!

# Adaptive Moment Estimation (Adam)

Combines Momentum and RMSProp

$$m^{k+1} = \beta_1 \cdot m^k + (1 - \beta_1)\nabla_\theta L(\theta^k)$$

First momentum: mean of gradients

$$v^{k+1} = \beta_2 \cdot v^k + (1 - \beta_2)[\nabla_\theta L(\theta^k) \circ \nabla_\theta L(\theta^k)]$$

Second momentum: variance of gradients

$$\theta^{k+1} = \theta^k - \alpha \cdot \frac{m^{k+1}}{\sqrt{v^{k+1}} + \epsilon}$$

# Adam

## Combines Momentum and RMSProp

$$m^{k+1} = \beta_1 \cdot m^k + (1 - \beta_1)\nabla_\theta L(\theta^k)$$

$m^{k+1}$ and $v^{k+1}$ are initialized with zero
-> bias towards zero

$$v^{k+1} = \beta_2 \cdot v^k + (1 - \beta_2)[\nabla_\theta L(\theta^k) \circ \nabla_\theta L(\theta^k)]$$
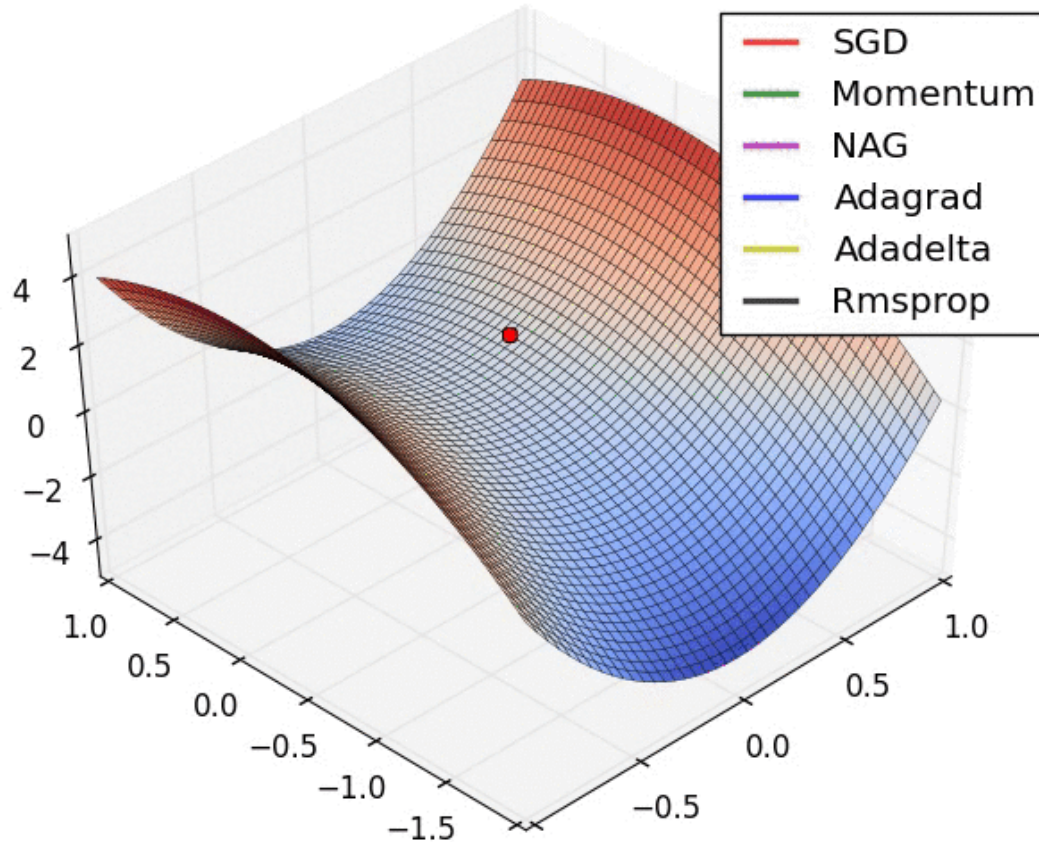
Typically, bias-corrected moment updates

$$\hat{m}^{k+1} = \frac{m^k}{1 - \beta_1}$$

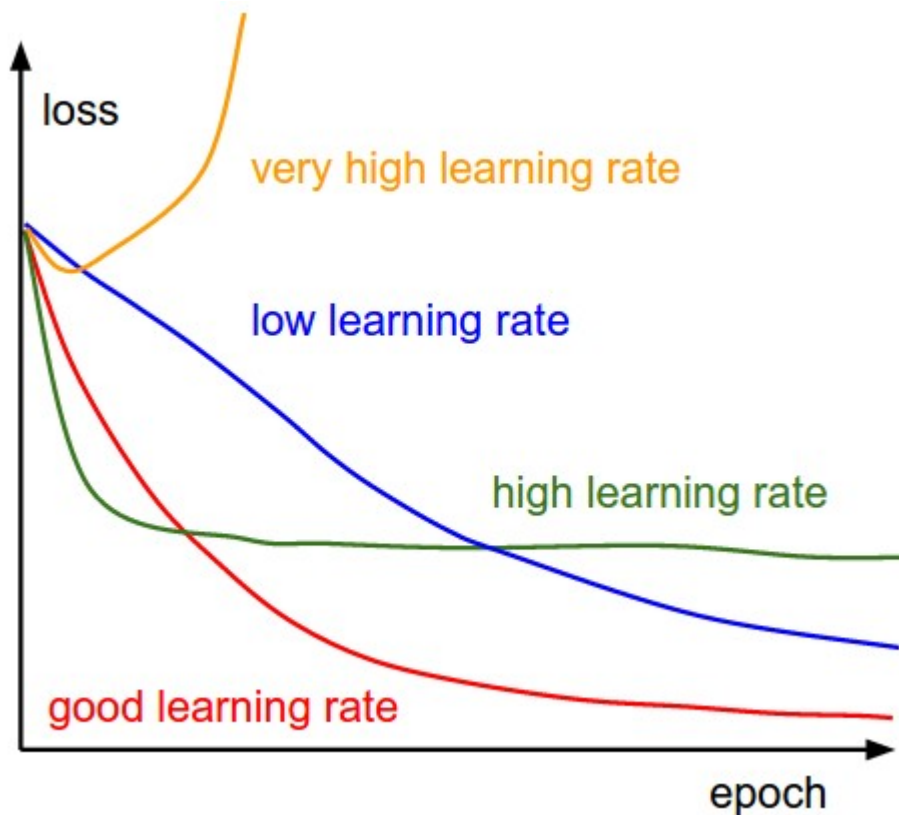$$\theta^{k+1} = \theta^k - \alpha \cdot \frac{\hat{m}^{k+1}}{\sqrt{\hat{v}^{k+1}} + \epsilon}$$

$$\hat{v}^{k+1} = \frac{v^k}{1 - \beta_2}$$

# Convergence

Prof. Leal-Taixé and Prof. Niessner

# Convergence

# Importance of Learning Rate

# Jacobian and Hessian

- Derivative $\quad \mathbf{f}: \mathbb{R} \to \mathbb{R} \qquad \dfrac{df(x)}{dx}$

- Gradient $\quad \mathbf{f}: \mathbb{R}^m \to \mathbb{R} \quad \nabla_{\mathbf{x}} f(\mathbf{x}) \quad \left( \dfrac{df(x)}{dx_1}, \dfrac{df(x)}{dx_2} \right)$

- Jacobian $\quad \mathbf{f}: \mathbb{R}^m \to \mathbb{R}^n \quad \mathbf{J} \in \mathbb{R}^{n \times m}$

- Hessian $\quad \mathbf{f}: \mathbb{R}^m \to \mathbb{R} \quad \mathbf{H} \in \mathbb{R}^{m \times m}$ SECOND DERIVATIVE

Prof. Leal-Taixé and Prof. Niessner

# Newton's method

- Approximate our function by a second-order Taylor series expansion

$$L(\boldsymbol{\theta}) \approx L(\boldsymbol{\theta}_0) + (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_0) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \mathbf{H}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)$$

First derivative

Second derivative (curvature)

https://en.wikipedia.org/wiki/Taylor_series

# Newton's method

- Differentiate and equate to zero

$$\boldsymbol{\theta}^* = \boldsymbol{\theta}_0 \boxed{- \mathbf{H}^{-1} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})} \quad \text{Update step}$$

We got rid of the learning rate!

$$\text{SGD} \quad \boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \epsilon \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_k, \mathbf{x}^i, \mathbf{y}^i)$$

# Newton's method

- Differentiate and equate to zero

$$\boldsymbol{\theta}^* = \boldsymbol{\theta}_0 \boxed{- \mathbf{H}^{-1} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})} \quad \text{Update step}$$

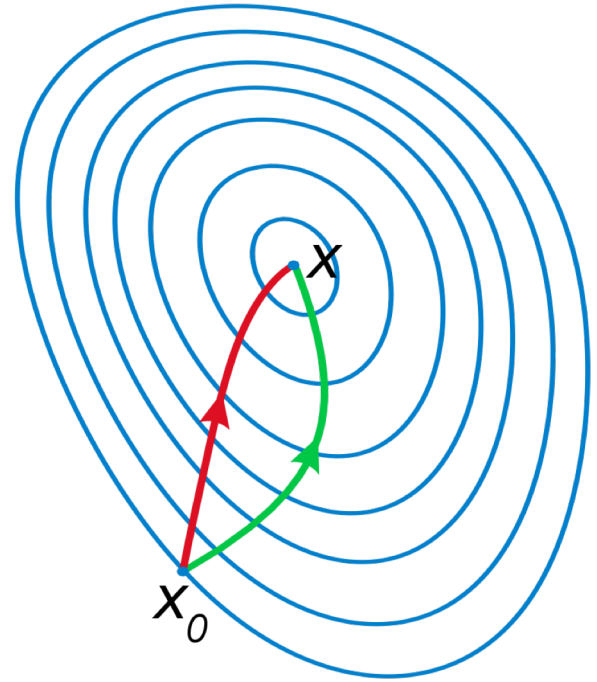| Parameters of a network (millions) | Number of elements in the Hessian | Computational complexity of 'inversion' per iteration |
|:---:|:---:|:---:|
| $k$ | $k^2$ | $\mathcal{O}(k^3)$ |

# Newton's method

- SGD (green)

- Newton's method exploits the curvature to take a more direct route

Image from Wikipedia

# Newton's method

$$J(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})$$

Can you apply Newton's method for linear regression? What do you get as a result?

# BFGS and L-BFGS

- Broyden-Fletcher-Goldfarb-Shanno algorithm
- Belongs to the family of quasi-Newton methods
- Have an approximation of the inverse of the Hessian

$$\boldsymbol{\theta}^* = \boldsymbol{\theta}_0 - \boxed{\mathbf{H}^{-1}} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

- BFGS $\quad \mathcal{O}(n^2)$
- Limited memory: L-BFGS $\quad \mathcal{O}(n)$

# Gauss-Newton

- $x_{k+1} = x_k - H_f(x_k)^{-1} \nabla f(x_k)$
  - 'true' 2nd derivatives are often hard to obtain (e.g., numerics)
  - $H_f \approx 2 J_F^T J_F$
- Gauss-Newton (GN):
$$x_{k+1} = x_k - [2 J_F(x_k)^T J_F(x_k)]^{-1} \nabla f(x_k)$$

- Solve linear system (again, inverting a matrix is unstable):
$$2\left(J_F(x_k)^T J_F(x_k)\right)\underbrace{(x_k - x_{k+1})}_{\text{Solve for delta vector}} = \nabla f(x_k)$$

# Levenberg

- Levenberg
  - "damped" version of Gauss-Newton:
  - $(J_F(x_k)^T J_F(x_k) + \lambda \cdot I) \cdot (x_k - x_{k+1}) = \nabla f(x_k)$

    **Tikhonov regularization**

  - The damping factor $\lambda$ is adjusted in each iteration ensuring:
  - $$f(x_k) > f(x_{k+1})$$
    - if inequation is not fulfilled increase $\lambda$
    - →Trust region

- →"Interpolation" between Gauss-Newton (small $\lambda$) and Gradient Descent (large $\lambda$)

# Levenberg-Marquardt

- Levenberg-Marquardt (LM)

$$(J_F(x_k)^T J_F(x_k) + \lambda \cdot diag(J_F(x_k)^T J_F(x_k))) \cdot (x_k - x_{k+1})$$
$$= \nabla f(x_k)$$

- Instead of a plain Gradient Descent for large $\lambda$, scale each component of the gradient according to the curvature.
  - Avoids slow convergence in components with a small gradient

# Which, what and when?

- Standard: **Adam**

- Fallback option: **SGD with momentum**

- **Newton, L-BFGS, GN, LM** only if you can do full batch updates (doesn't work well for minibatches!!)

This practically never happens for DL
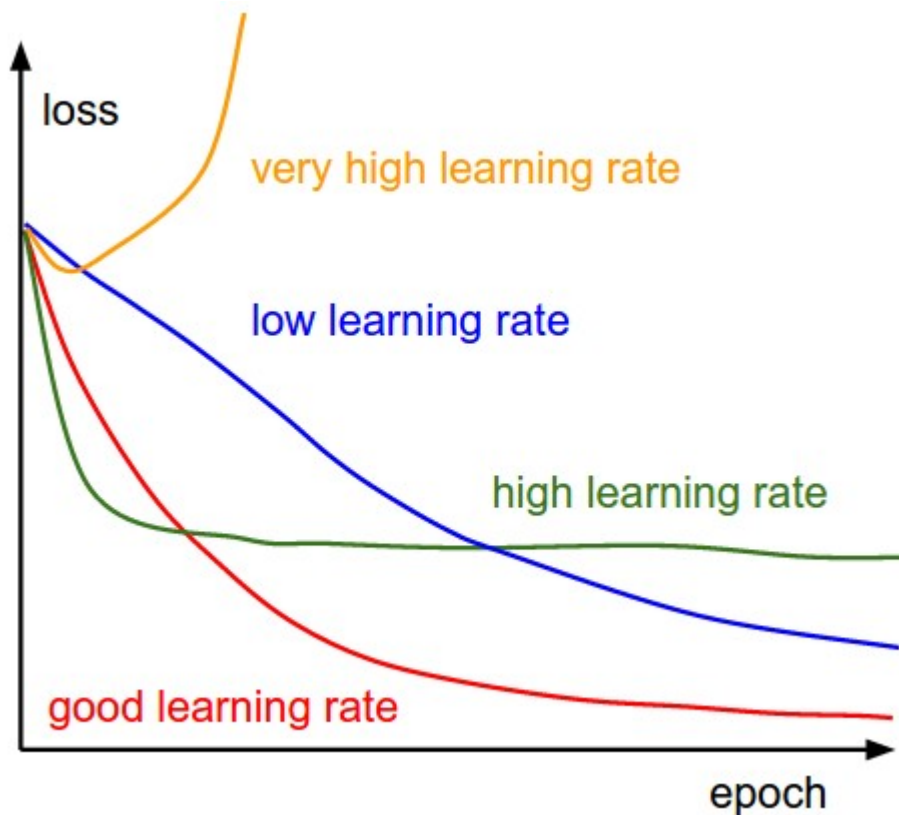Theoretically, it would be nice though due to fast convergence

# General Optimization

- Linear Systems (Ax = b)
  - LU, QR, Cholesky, Jacobi, Gauss-Seidel, CG, PCG, etc.
- Non-linear (gradient-based)
  - Newton, Gauss-Newton, LM, (L)BFGS     <- second order
  - Gradient Descent, SGD                  <- first order

- Others:
  - Genetic algorithms, MCMC, Metropolis-Hastings, etc.
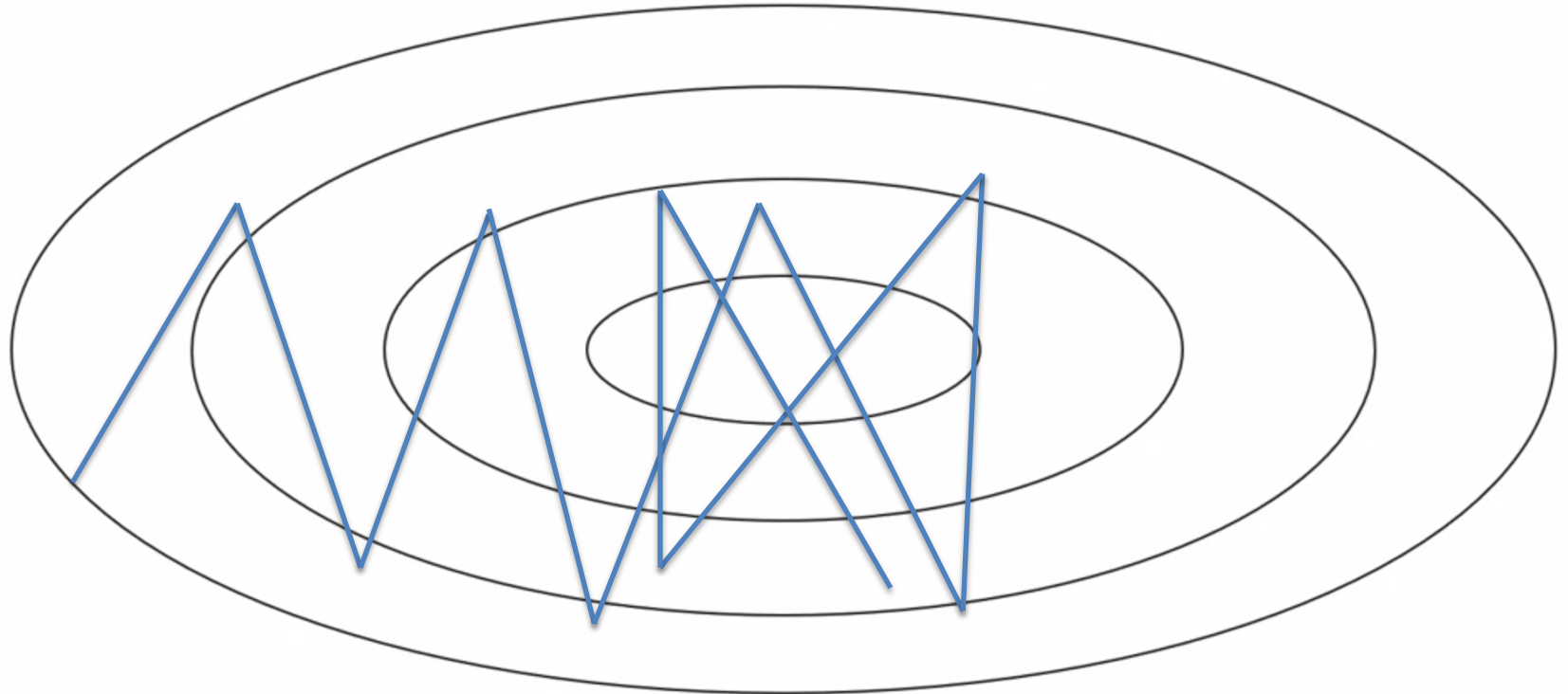  - Constrained and convex solvers (Langrage, ADMM, Primal-Dual, etc.)

# Please Remember!

- Think about your problem and optimization at hand

- SGD is specifically designed for minibatch

- When you can, use $2^{nd}$ order method -> it's just faster

- GD or SGD is **not** a way to solve a linear system!

# Importance of Learning Rate

# Learning Rate



Need high learning rate when far away    Need low learning rate when close

# Learning Rate Decay

- $\alpha = \dfrac{1}{1 + decayrate \cdot epoch} \cdot \alpha_0$
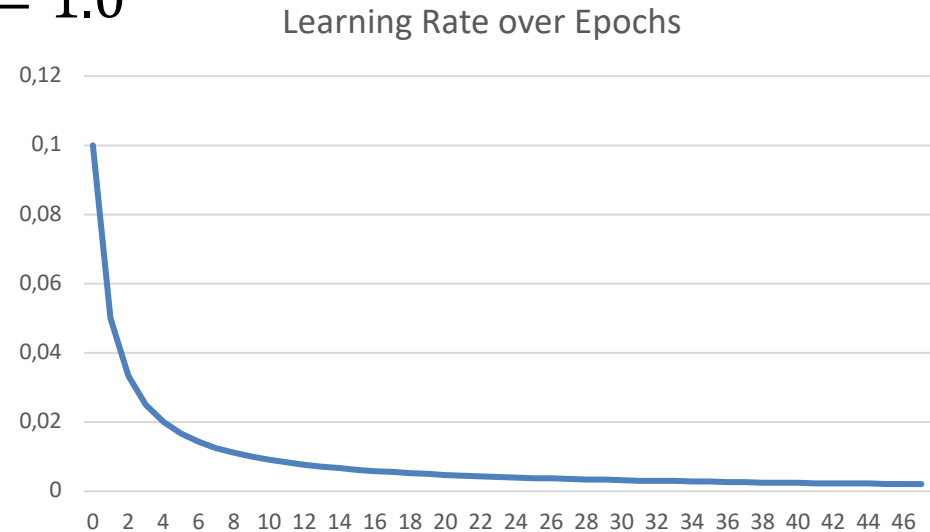
  – E.g., $\alpha_0 = 0.1$, $decayrate = 1.0$

  – > Epoch 0:  0.1
  – > Epoch 1:  0.05
  – > Epoch 2:  0.033
  – > Epoch 3:  0.025

    ...



Learning Rate over Epochs

# Learning Rate Decay

Many options:

- Step decay $\alpha = \alpha - t \cdot \alpha$ (only every n steps)
  - T is decay rate (often 0.5)

- Exponential decay $\alpha = t^{epoch} \cdot \alpha_0$
  - t is decay rate (t < 1.0)

- $\alpha = \dfrac{t}{\sqrt{epoch}} \cdot a_0$
  - t is decay rate
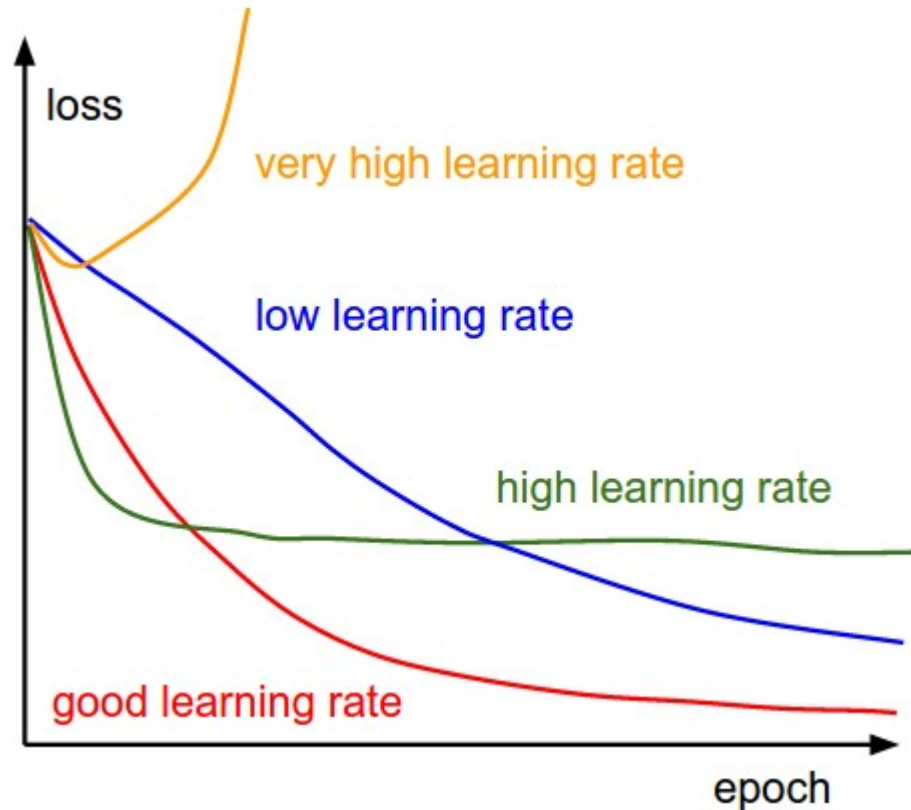
- Etc.

# Training Schedule

Manually specify learning rate for entire training process

- Manually set learning rate every n-epochs
- How?
  - Trial and error (the hard way)
  - Some experience (only generalizes to some degree)

Consider: #epochs, training set size, network size, etc.

# Learning Rate: Implications

- What if too high?

- What if too low?

# Training

- Given ground dataset with ground lables
  - $\{x_i, y_i\}$
    - For instance $x_i$-th training image, with label $y_i$
    - Often $\dim(x) \gg \dim(y)$ (e.g., for classification)
    - $i$ is often in the 100-thousands or millions
  - Take network $f$ and its parameters $w, b$

  - Use SGD (or variation) to find optimal parameters $w, b$
    - Gradients from backprop

# Learning

- Learning means generalization to unknown dataset
  - (so far no 'real' learning)
  - I.e., train on known dataset -> test with optimized parameters on unknown dataset


- Basically, we hope that based on the train set, the optimized parameters will give similar results on different data (i.e., test data)
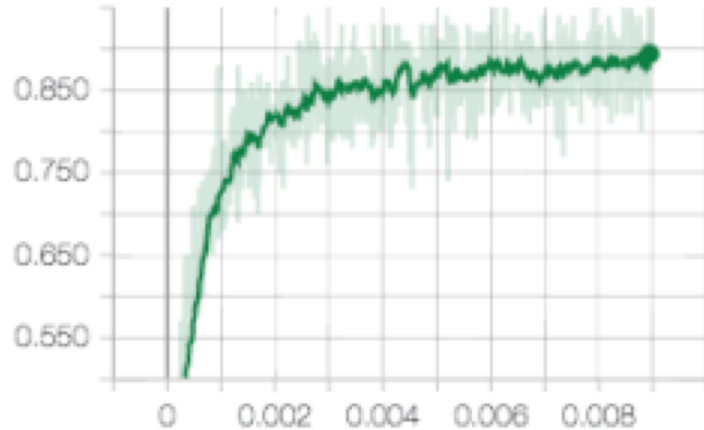
# Learning

- Training set ('*train*'):
  - Use for training your neural network
- Validation set ('*val*'):
  - Hyperparameter optimization
  - Check generalization progress
- Test set ('*test*'):
  - Only for the very end
  - NEVER TOUCH DURING DEVELOPMENT OR TRAINING
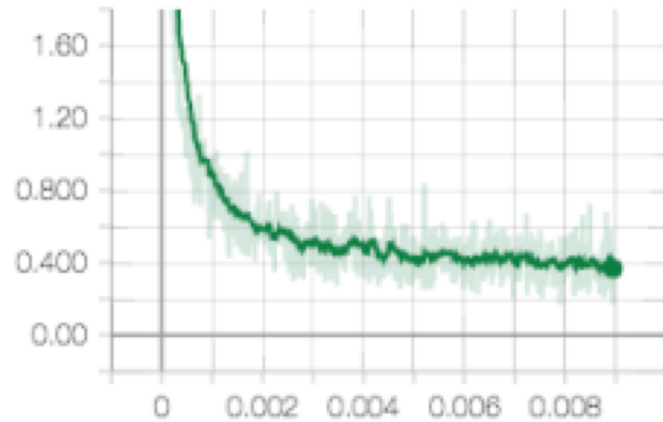
# Learning

- Typical splits
  - Train (60%), Val (20%), Test (20%)
  - Train (80%), Val (10%), Test (10%)

- During training:
  - Train error comes from average mini-batch error
  - Typically take subset of validation every n iterations
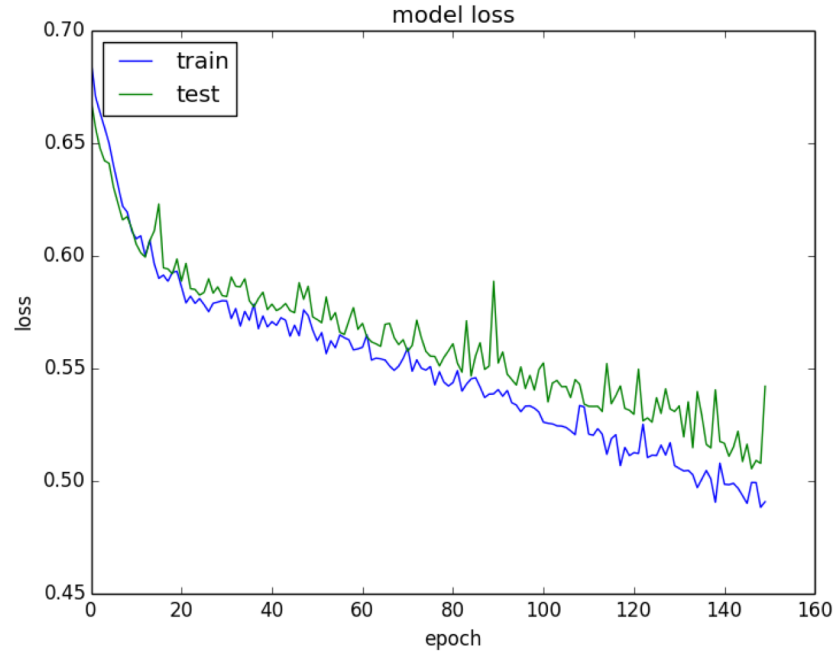
# Learning

- Training graph
  - Accuracy                                    - Loss



(EMA smoothing)
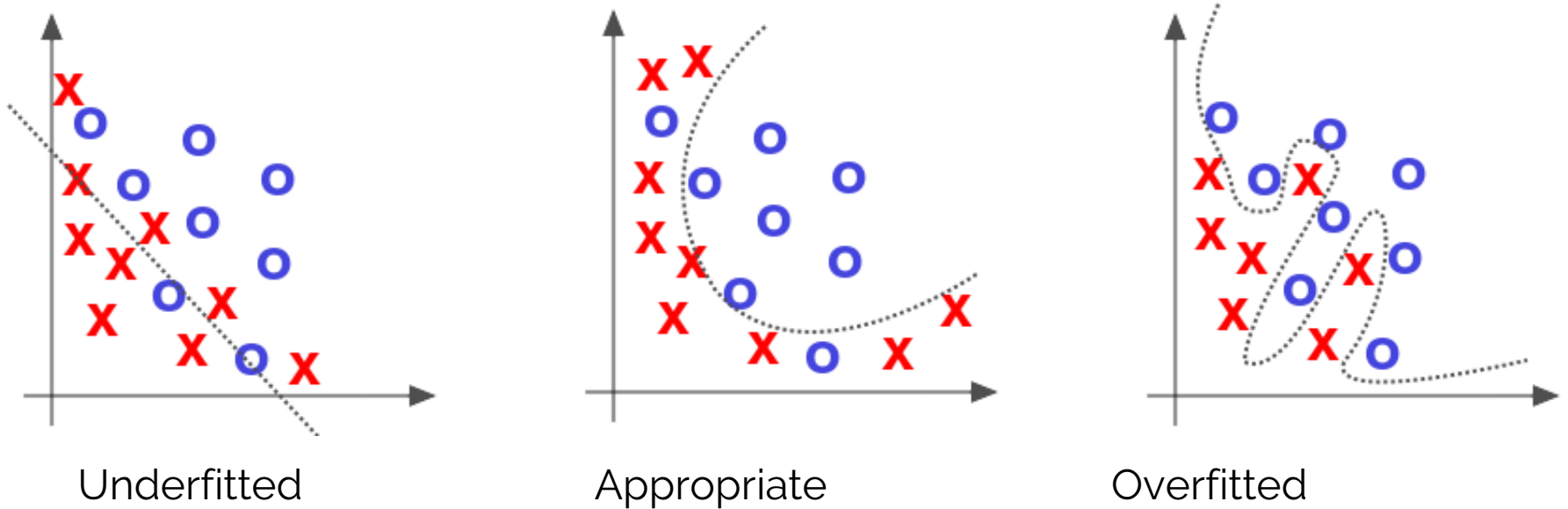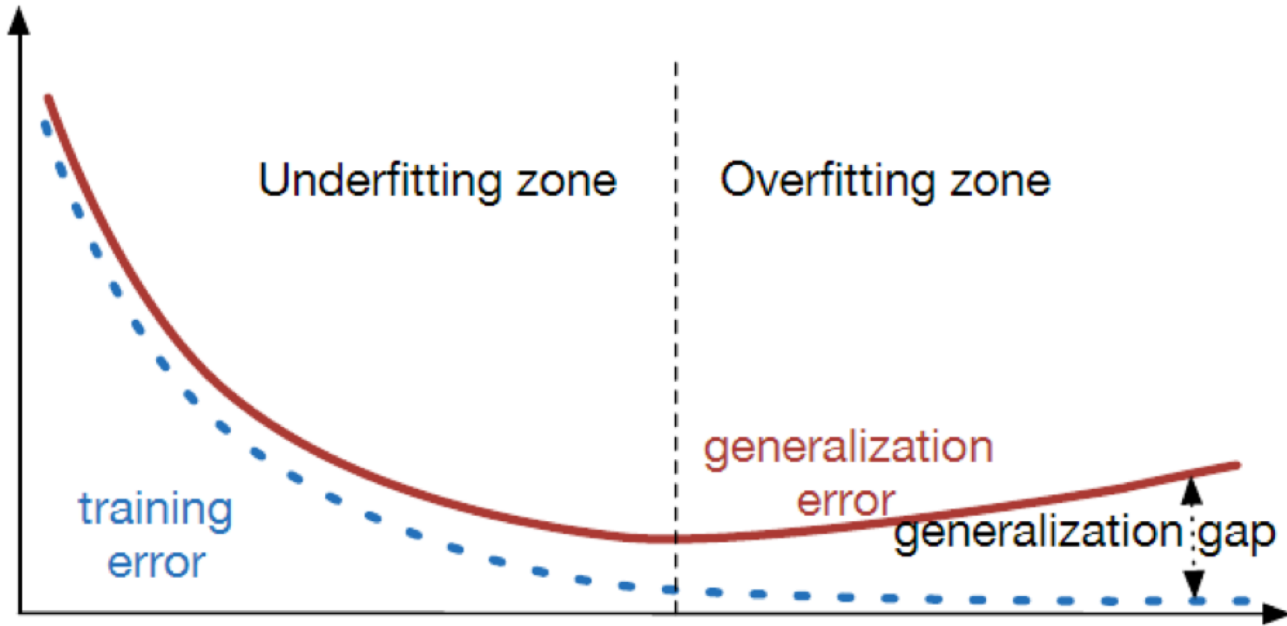
# Learning

- Validation graph

# Over- and Underfitting



Underfitted        Appropriate        Overfitted

Figure extracted from Deep Learning by Adam Gibson, Josh Patterson, O'Reily Media Inc., 2017

# Over- and Underfitting



Underfitting zone | Overfitting zone

generalization
error

training
error

generalization gap

Source: http://srdas.github.io/DLBook/ImprovingModelGeneralization.html

# Hyperparameters

- Network architecture (e.g., num layers, #weights)

- Number of iterations

- Learning rate(s)  (i.e., solver parameters, decay, etc.)

- Regularization (more later next lecture)

- Batch size

- ...

- Overall: learning setup + optimization = hyerparameter

# Hyperparameter Tuning

- Methods:
  - Manual search: most common ☺
  - Grid search (structured, for 'real' applications)

    Define ranges for all parameters spaces and select points (usually pseudo-uniformly distributed). Iterate over all  possible configurations

  - Random search:

    Like grid search but one picks points at random in the predefined ranges

# Simple Grid Search Example

```
learning_rates = [1e-2, 1e-3, 1e-4, 1e-5]
regularization_strengths = [1e2, 1e3, 1e4, 1e5]
num_iters = [500, 1000, 1500]
best_val = 0

for learning_rate in learning_rates:
  for reg in regularization_strengths:
    for iterations in num_iters:
        model = train_model(learning_rate, reg., iterations)
        validation_accuracy = evaluate(model)
        if validation_accuracy > best_val:
            best_val = validation_accuracy
            best_model = model
```
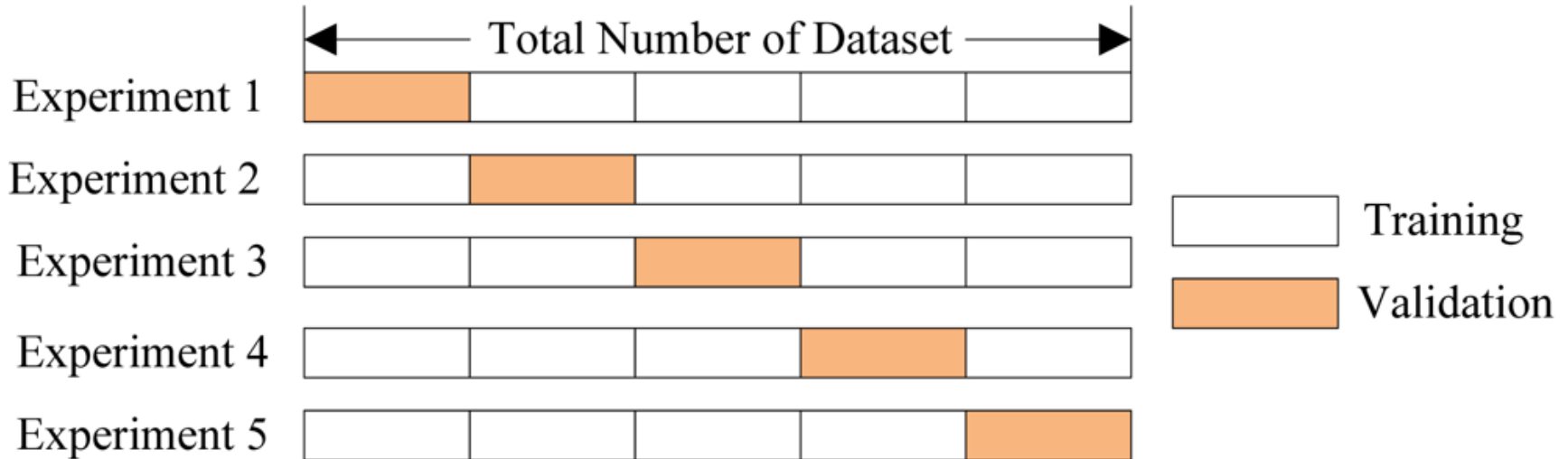
# Cross Validation

- Example: k=5



Figure extracted from cs231n

# Cross Validation

- Used when data set is extremely small and/or our method of choice has low training times

- Partition data into k subsets, train on k-1 and evaluate performance on the remaining subset

- To reduce variability: perform on different partitions and average results
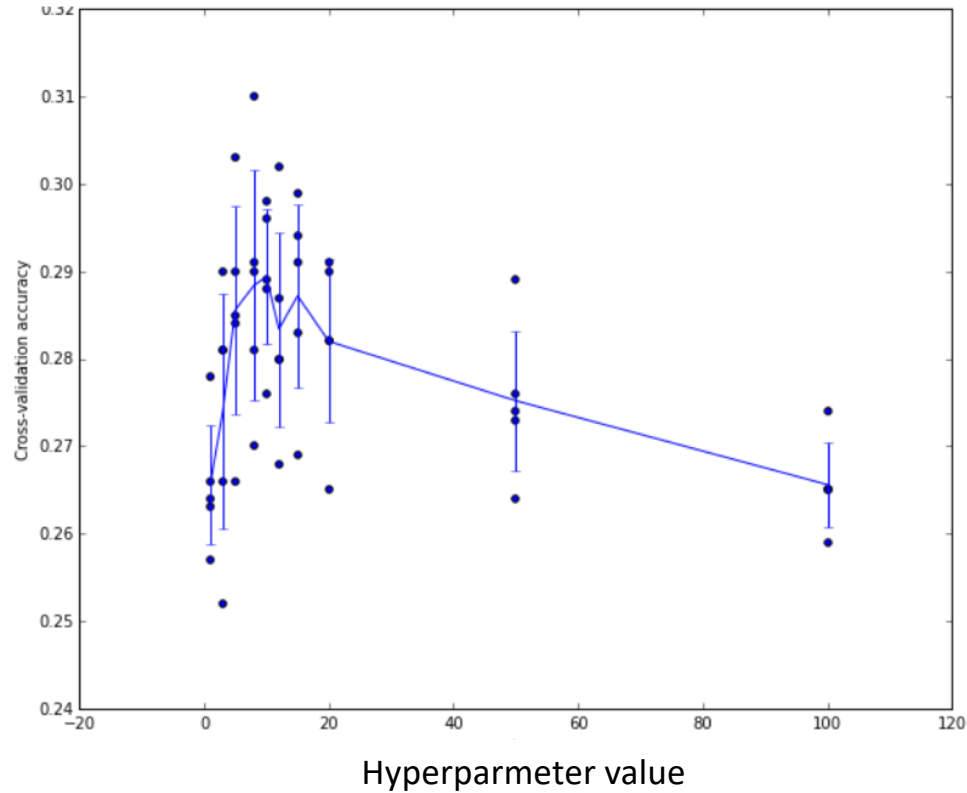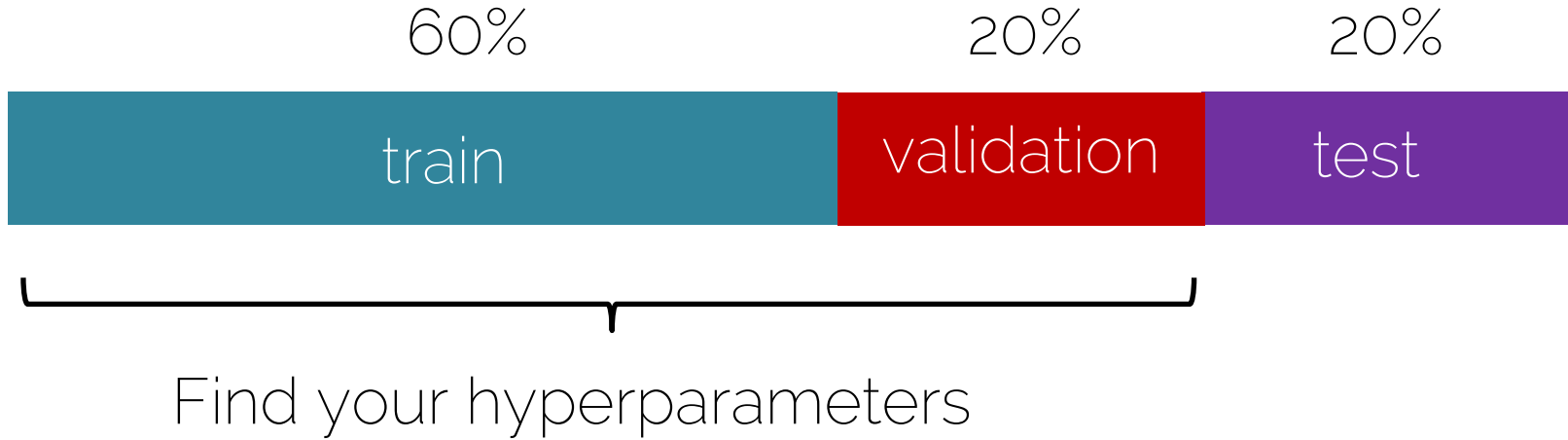
# Cross Validation

Results for k=5



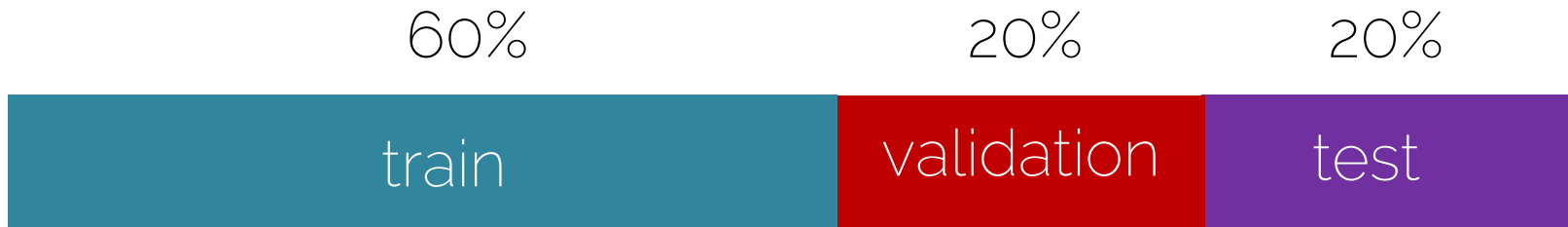Figure extracted from cs231n

# Basic recipe for machine learning

# Basic recipe for machine learning

- Split your data



| 60% | 20% | 20% |
| --- | --- | --- |
| train | validation | test |

Find your hyperparameters

# Basic recipe for machine learning

- Split your data

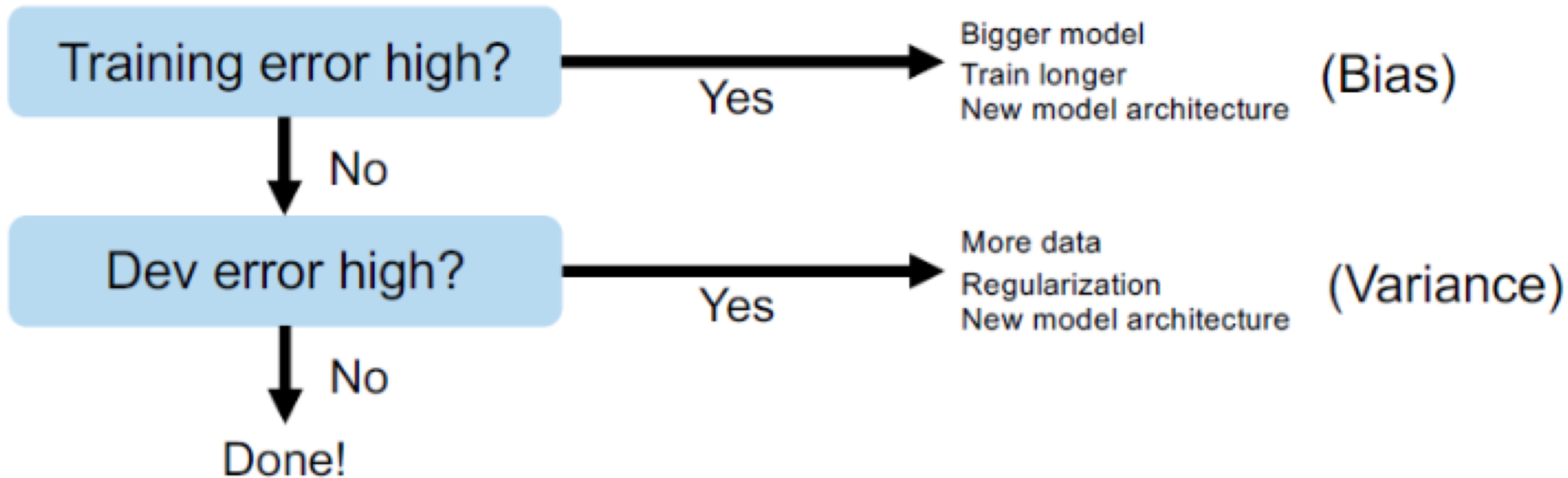| 60% | 20% | 20% |
|---|---|---|
| train | validation | test |

Human level error ...... 1%

Training set error ...... 5%

Val/test set error ...... 8%

*Bias* (or underfitting)

*Variance* (overfitting)

# Basic recipe for machine learning



| Training error high? | Yes → | Bigger model<br>Train longer<br>New model architecture | (Bias) |

No ↓

| Dev error high? | Yes → | More data<br>Regularization<br>New model architecture | (Variance) |

No ↓

Done!

More on

# Next lecture

- This week:
  - Maybe exercise session (see upcoming moodle announcement, NIPS deadline)


- Next lecture on May 21$^{st}$:
  - Training Neural Networks

Prof. Leal-Taixé and Prof. Niessner

# See you next week!