

# Lecture 15: Optimization (Under/over)fitting

# Administrative

- HW3 due Wednesday, March 4 11:59pm
- TAs will not be checking Piazza over Spring Break. You are **strongly encouraged** to finish the assignment by Friday, February 25

# Last Time: Regularized Least Squares

Add **regularization** to objective that prefers some solutions:

Before:  $\arg \min_w \|y - Xw\|_2^2 \longrightarrow \text{Loss}$

After:  $\arg \min_w \|y - Xw\|_2^2 + \lambda \|w\|_2^2$

Loss      Trade-off      Regularization



Want model “smaller”: pay a penalty for  $w$  with big norm

Intuitive Objective: accurate model (low loss) but not too complex (low regularization).  $\lambda$  controls how much of each.

# Last Time: Nearest Neighbor

Known Images

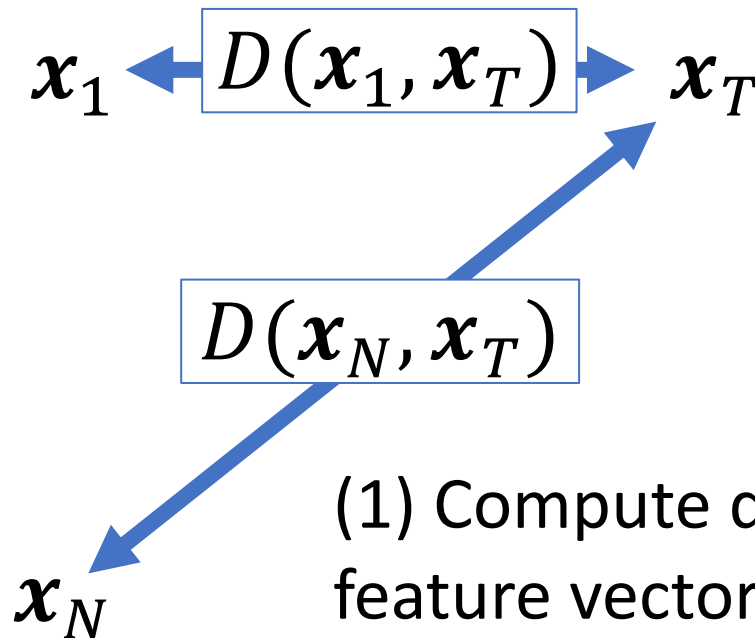
Labels



...



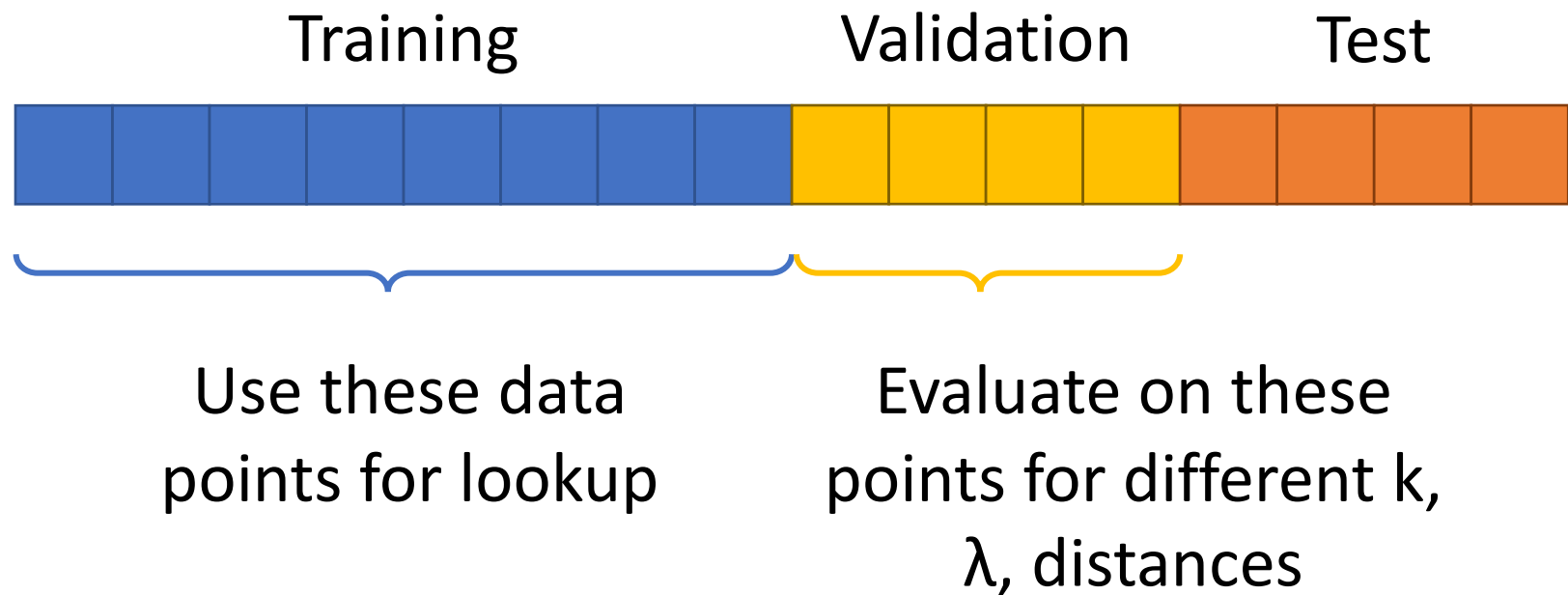
Test Image



- (1) Compute distance between feature vectors
- (2) find nearest
- (3) use label.

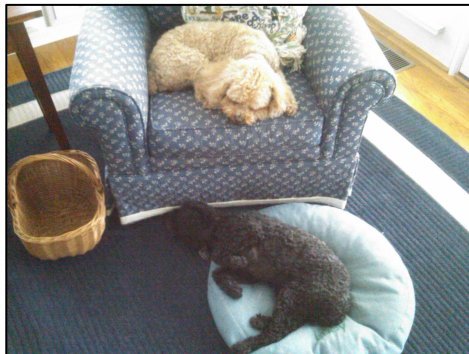
# Last Time: Choosing Hyperparameters

What distance? What value for  $k$  /  $\lambda$ ?

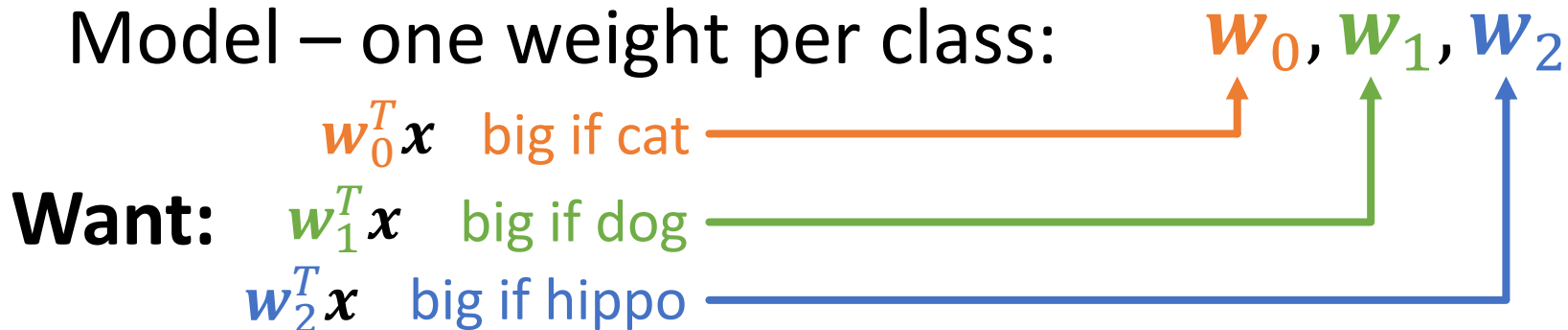


# Last Time: Linear Classifiers

## Example Setup: 3 classes



Model – one weight per class:



Stack together:  $W_{3 \times F}$  where  $x$  is in  $\mathbb{R}^F$

# Last Time: Linear Classifiers

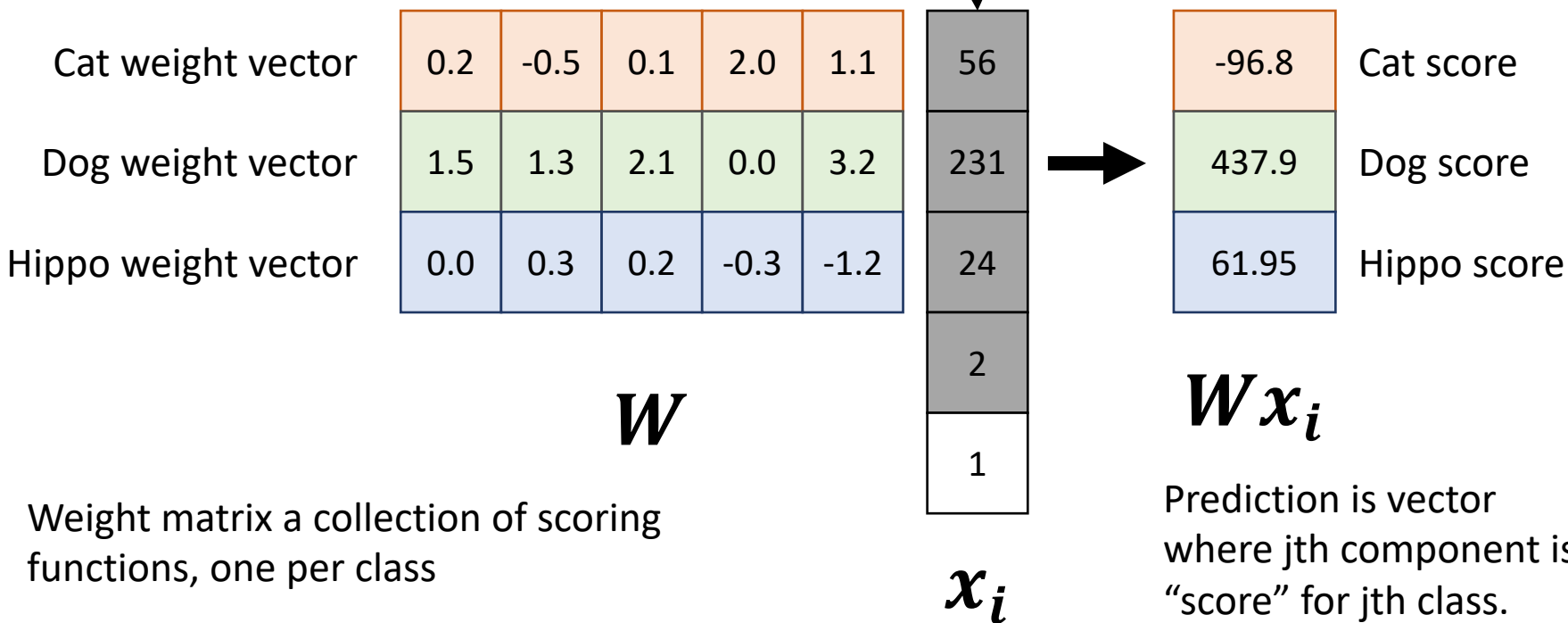
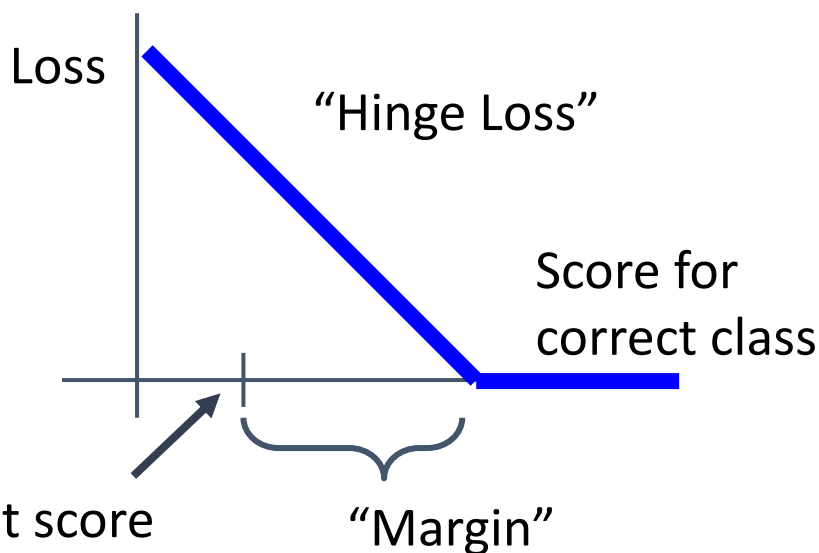


Diagram by: Karpathy, Fei-Fei

# Last Time: Multiclass SVM Loss

“The score of the correct class should be higher than all the other scores”



Given an example  $(x_i, y_i)$   
( $x_i$  is image,  $y_i$  is label)

Let  $s = f(x_i, W)$  be scores

Then the SVM loss has the form:

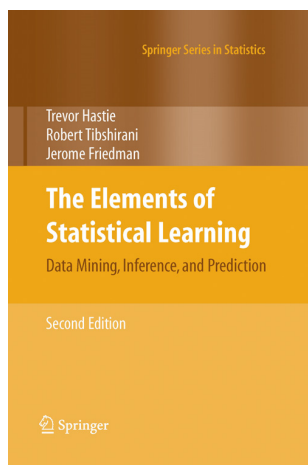
$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$



# Last Time: Multiclass SVM Loss

SVM = Support Vector Machine

Lots of great theory as to why this is a sensible thing to do. See



Useful book (Free too!):

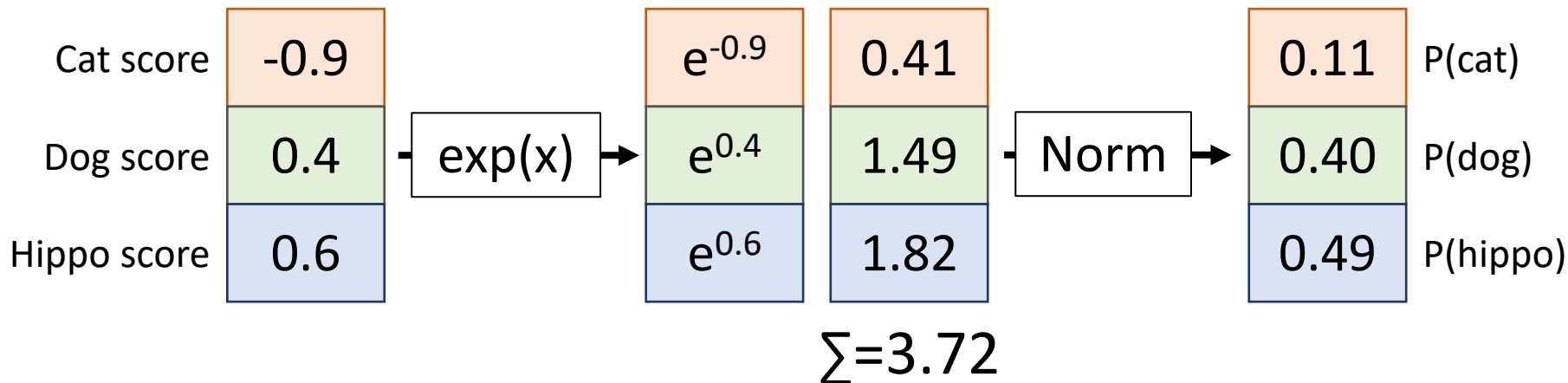
The Elements of Statistical Learning

Hastie, Tibshirani, Friedman

<https://web.stanford.edu/~hastie/ElemStatLearn/>

# Last Time: Cross-Entropy Loss

Converting Scores to “Probability Distribution”



Generally P(class j):

$$\frac{\exp((Wx)_j)}{\sum_k \exp((Wx)_k)}$$

*Called softmax function*

Loss is  $-\log(P(\text{correct class}))$

$$L_i = -\log \frac{\exp(s_{y_i})}{\sum_j \exp(s_j)}$$

# Today: Optimization

Goal: find the  $\mathbf{w}$  minimizing some loss function  $L$ .

$$\arg \min_{\mathbf{w} \in \mathbb{R}^N} L(\mathbf{w})$$

Works for lots of different  $L$ s:

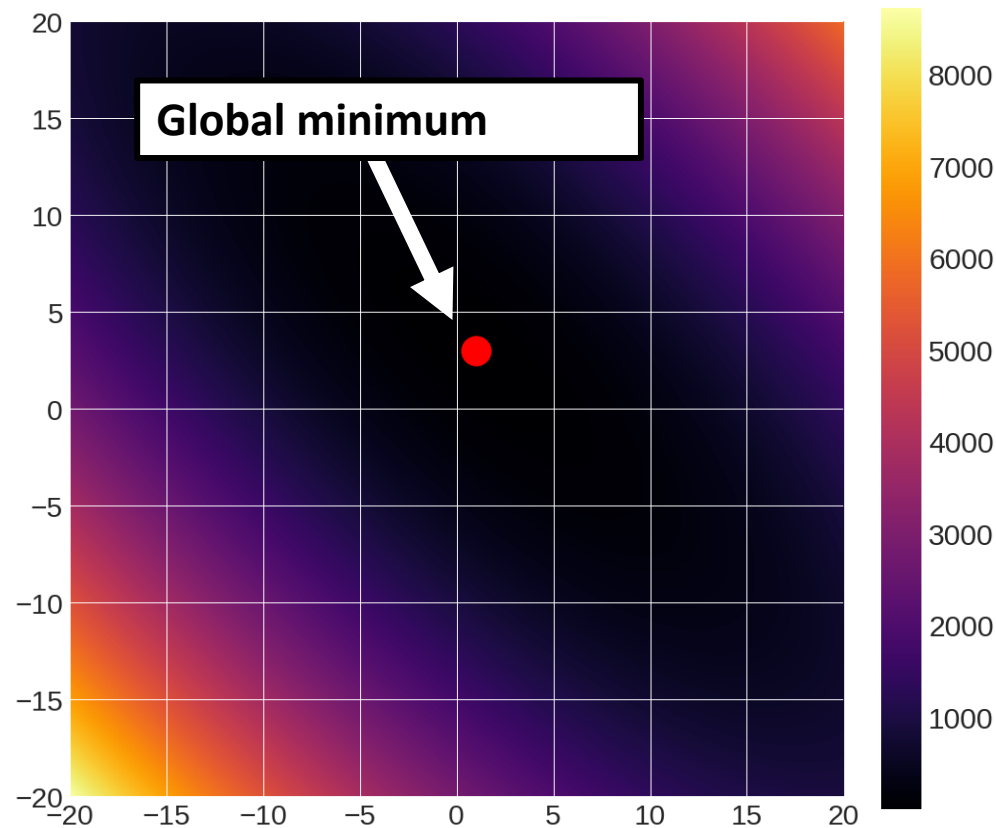
$$L(\mathbf{W}) = \lambda \|\mathbf{W}\|_2^2 + \sum_{i=1}^n -\log \left( \frac{\exp((\mathbf{W}\mathbf{x})_{y_i})}{\sum_k \exp((\mathbf{W}\mathbf{x})_k)} \right)$$

$$L(\mathbf{w}) = \lambda \|\mathbf{w}\|_2^2 + \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

$$L(\mathbf{w}) = C \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)$$

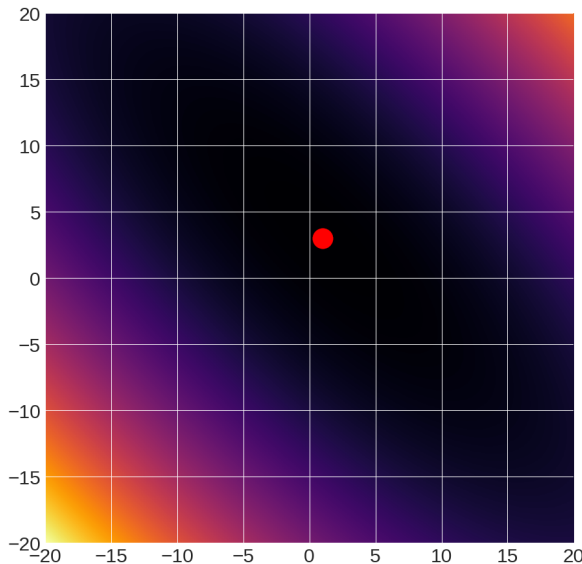
# Sample Function to Optimize

$$f(x,y) = (x+2y-7)^2 + (2x+y-5)^2$$



**Warning:** This is 2D, intuition may not generalize to high dimension

# Optimization: A Caveat



- Each point in the picture is a function evaluation
- Here it takes microseconds – so we can easily see the answer
- Functions we want to optimize may take hours to evaluate



[This image](#) is [CC0 1.0](#) public domain

[Walking man image](#) is [CC0 1.0](#) public domain



This image is [CC0 1.0](https://creativecommons.org/licenses/by/4.0/) public domain

Walking man image is [CC0 1.0](https://creativecommons.org/licenses/by/4.0/) public domain

# Idea #1A: Grid Search

#systematically try things

best, bestScore = None, Inf

for dim1Value in dim1Values:

....

for dimNValue in dimNValues:

$\mathbf{w} = [\text{dim1Value}, \dots, \text{dimNValue}]$

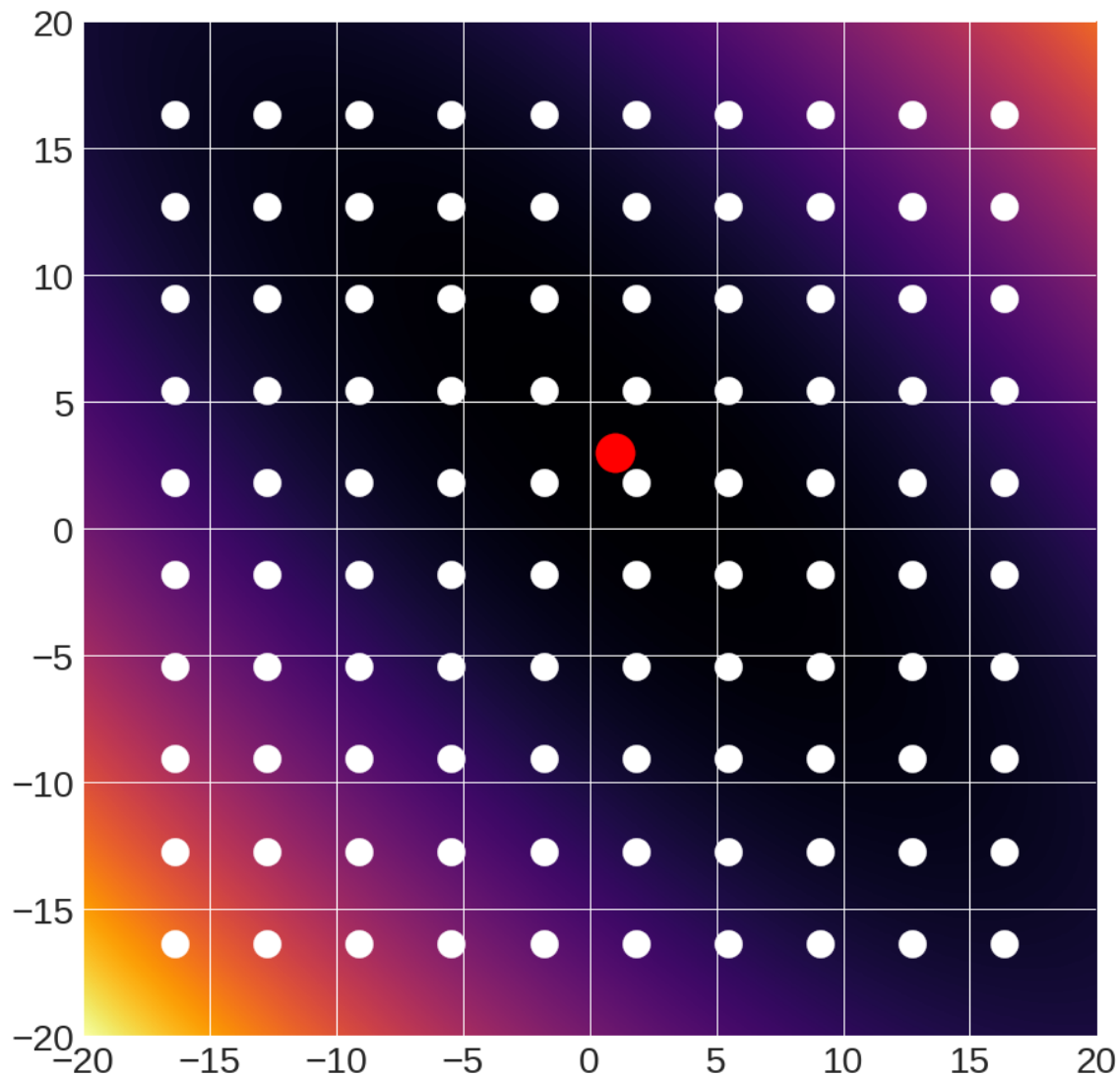
if  $L(\mathbf{w}) < \text{bestScore}$ :

best, bestScore =  $\mathbf{w}$ ,  $L(\mathbf{w})$

return best



# Idea #1A: Grid Search



# Idea #1A: Grid Search

## Pros:

1. Super simple
2. Only requires being able to evaluate model

## Cons:

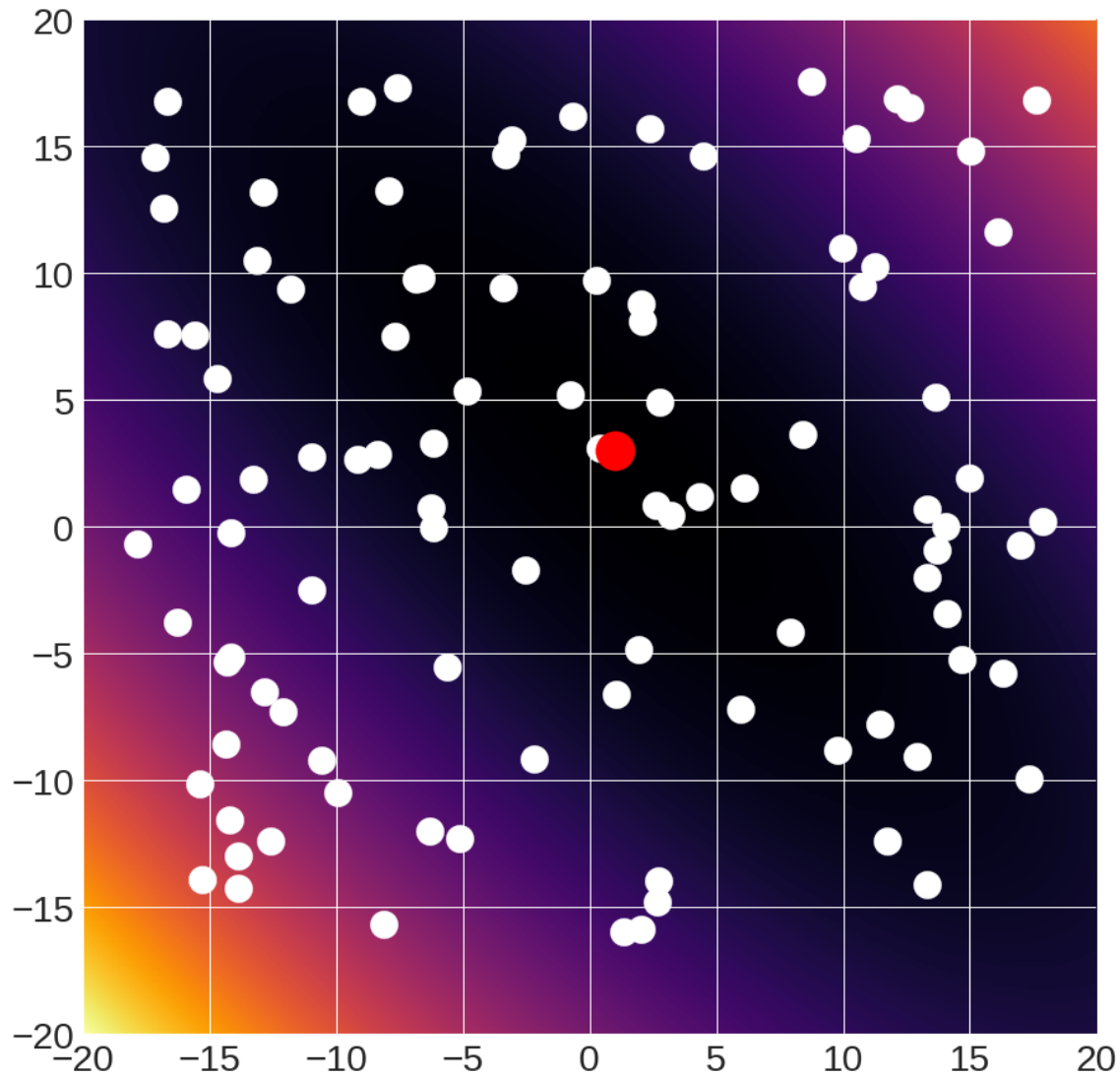
1. Scales horribly to high dimensional spaces

Complexity:  $\text{samplesPerDim}^{\text{numberOfDims}}$

# Option #1B: Random Search

```
#Do random stuff RANSAC Style
best, bestScore = None, Inf
for iter in range(numIters):
    w = random(N,1) #sample
    score =  $L(\mathbf{w})$  #evaluate
    if score < bestScore:
        best, bestScore = w, score
return best
```

# Option #1B: Random Search



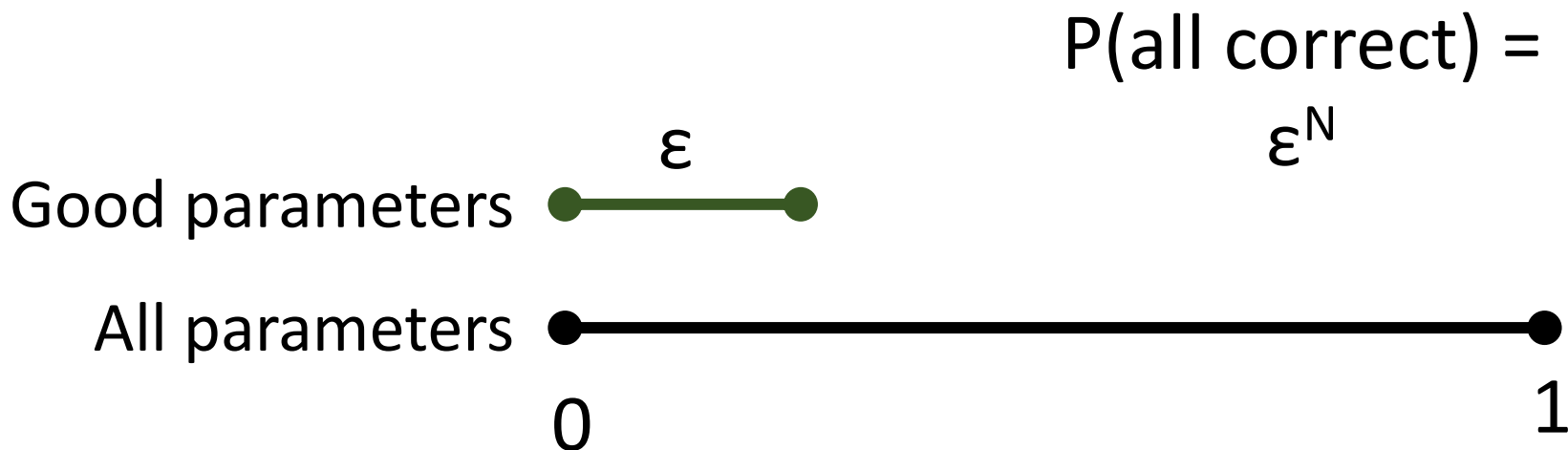
# Option #1B: Random Search

## Pros:

1. Super simple
2. Only requires being able to sample model and evaluate it

## Cons:

1. Slow –throwing darts at high dimensional dart board
2. Might miss something



# When To Use Options 1A / 1B?

Use these when

- Number of dimensions small, space bounded
- Objective is impossible to analyze (e.g., test accuracy if we use this distance function)

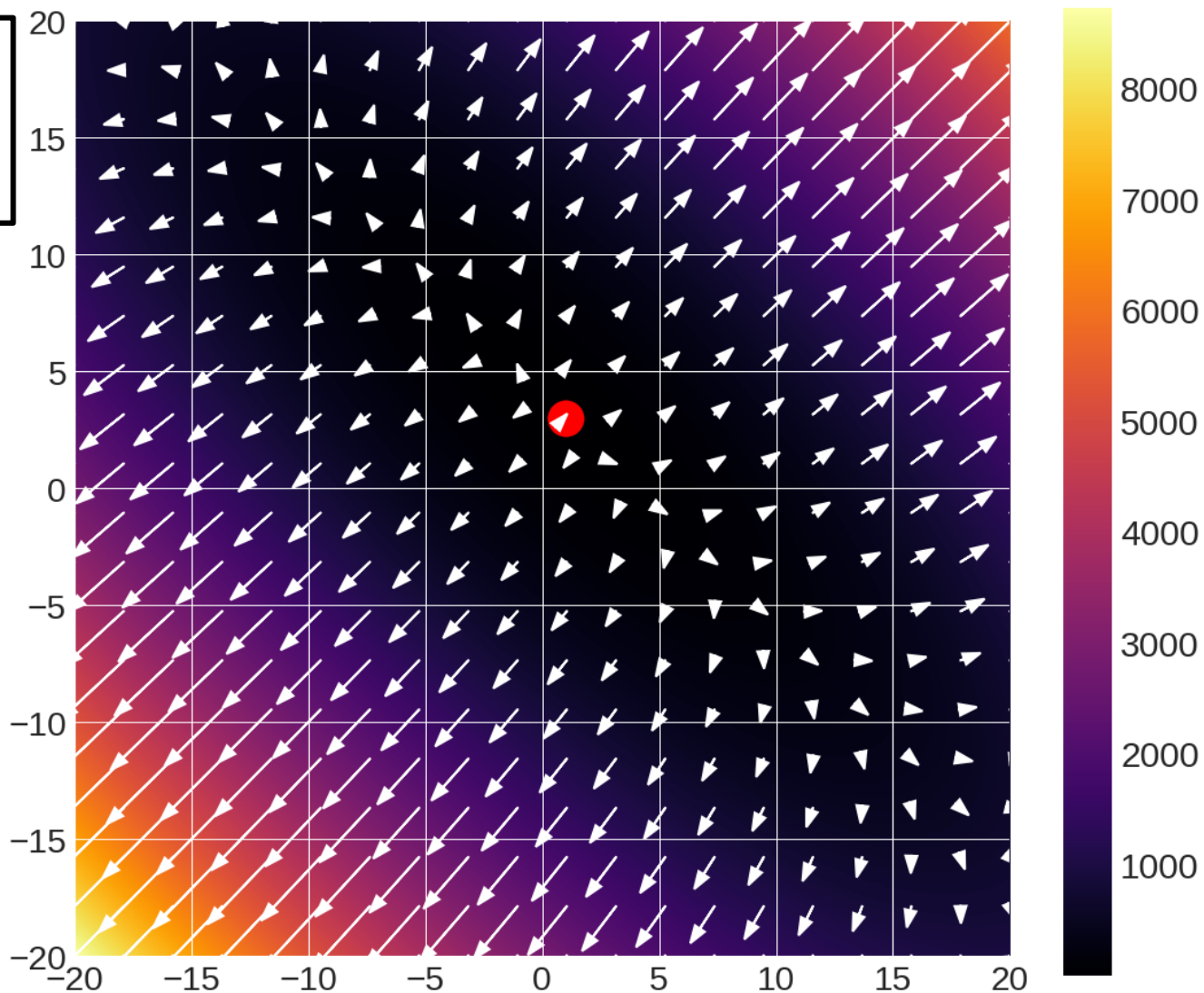
Random search is arguably more effective; grid search makes it easy to systematically test something (people love certainty)

# Idea #2: Follow the slope



# Idea #2: Follow the slope

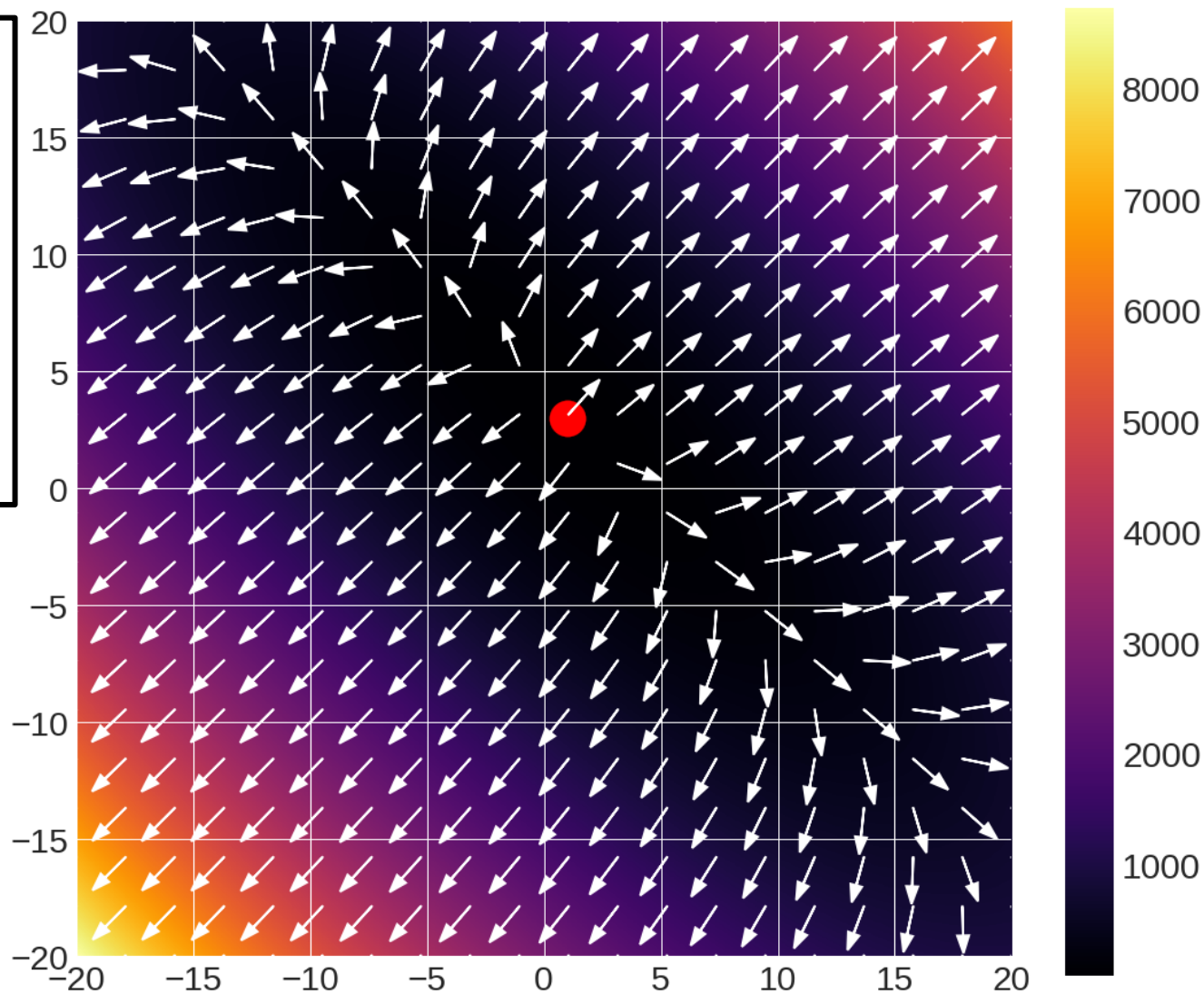
Arrows:  
gradient





# Idea #2: Follow the slope

Arrows:  
**gradient  
direction**  
(scaled to unit  
length)



# Idea #2: Follow the slope

Want:  $\arg \min_{\mathbf{w}} L(\mathbf{w})$

What's the geometric interpretation of:  $\nabla_{\mathbf{w}} L(\mathbf{w}) = \begin{bmatrix} \partial L / \partial \mathbf{x}_1 \\ \vdots \\ \partial L / \partial \mathbf{x}_N \end{bmatrix}$

Which is bigger (for small  $\alpha$ )?

$$L(\mathbf{w}) \begin{array}{l} \leq? \\ >? \end{array} L(\mathbf{w} + \alpha \nabla_{\mathbf{w}} L(\mathbf{w}))$$

# Idea #2: Follow the slope

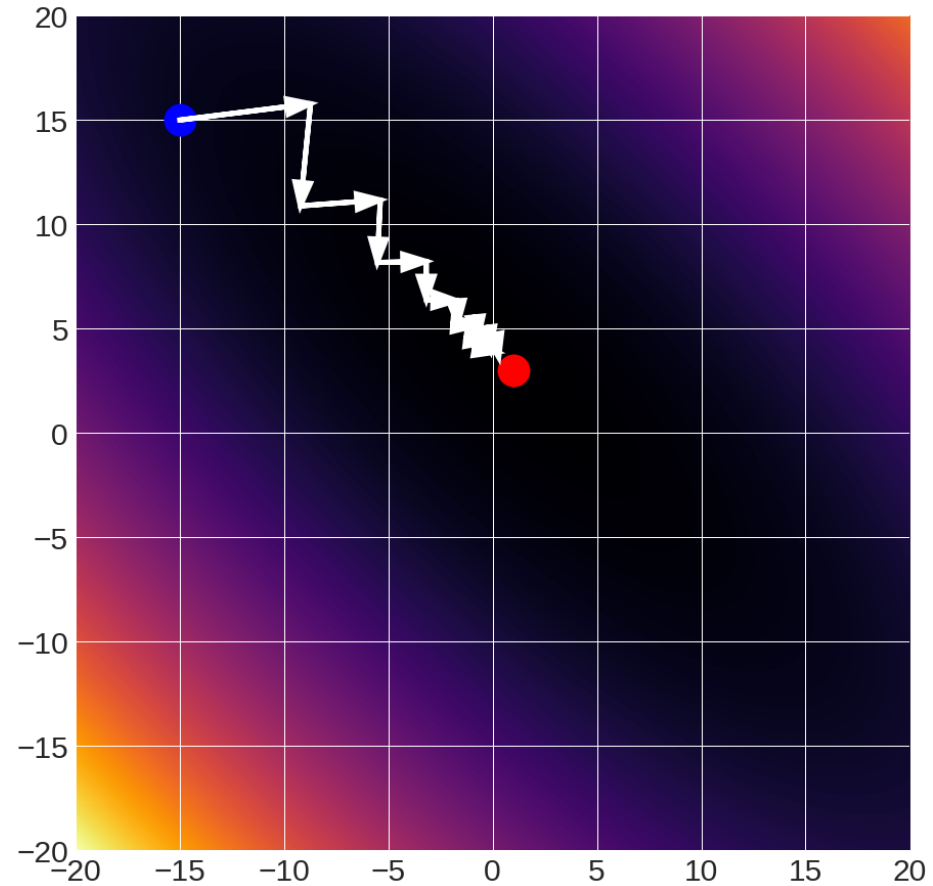
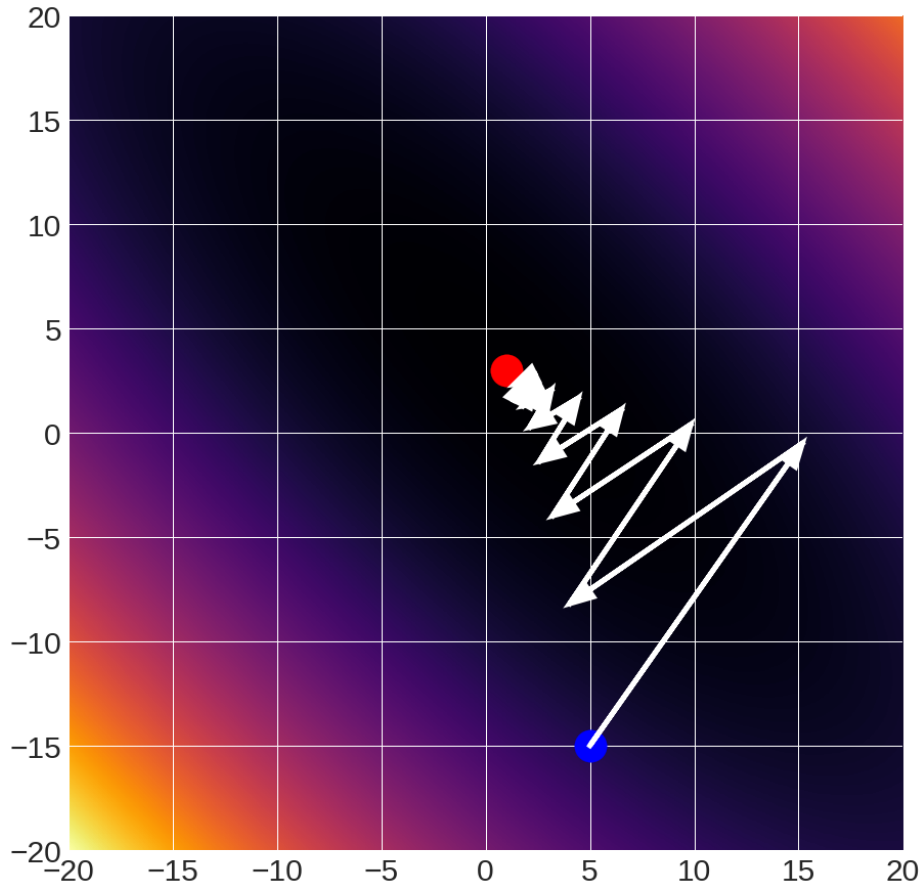
**Method:** at each step, move in direction of negative gradient

```
w0 = initialize() #initialize  
for iter in range(numIters):  
    g =  $\nabla_{\mathbf{w}}L(\mathbf{w})$  # eval gradient  
    w = w + -stepsize(iter)*g # update w  
return w
```

# Gradient Descent

Given starting point (blue)

$$w_{i+1} = w_i + -9.8 \times 10^{-2} \times \text{gradient}$$



# Computing Gradients: Numeric

How Do You Compute The Gradient?

Numerical Method:

$$\nabla_{\mathbf{w}} L(\mathbf{w}) = \begin{bmatrix} \frac{\partial L(\mathbf{w})}{\partial w_1} \\ \vdots \\ \frac{\partial L(\mathbf{w})}{\partial w_n} \end{bmatrix}$$

**How do you compute this?**

$$\frac{\partial f(x)}{\partial x} = \lim_{\epsilon \rightarrow 0} \frac{f(x + \epsilon) - f(x)}{\epsilon}$$

In practice, use:

$$\frac{f(x + \epsilon) - f(x - \epsilon)}{2\epsilon}$$

# Computing Gradients: Numeric

How Do You Compute The Gradient?

Numerical Method:

$$\nabla_{\mathbf{w}} L(\mathbf{w}) = \begin{bmatrix} \frac{\partial L(\mathbf{w})}{\partial x_1} \\ \vdots \\ \frac{\partial L(\mathbf{w})}{\partial x_n} \end{bmatrix}$$

$$\text{Use: } \frac{f(x + \epsilon) - f(x - \epsilon)}{2\epsilon}$$

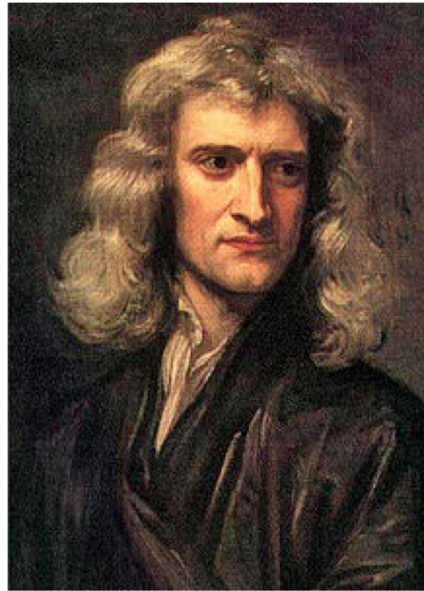
**How many function evaluations per dimension?**

# Computing Gradients: Analytic

How Do You Compute The Gradient?

Better Idea: Use Calculus!

$$\nabla_{\mathbf{w}} L(\mathbf{w}) = \begin{bmatrix} \frac{\partial L(\mathbf{w})}{\partial x_1} \\ \vdots \\ \frac{\partial L(\mathbf{w})}{\partial x_n} \end{bmatrix}$$



[This image](#) is in the public domain



[This image](#) is in the public domain

# Computing Gradients: Analytic

$$L(\mathbf{w}) = \lambda \|\mathbf{w}\|_2^2 + \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

$$\begin{array}{ccc} \downarrow & \frac{\partial}{\partial \mathbf{w}} & \downarrow \end{array}$$

$$\nabla_{\mathbf{w}} L(\mathbf{w}) = 2\lambda \mathbf{w} + \sum_{i=1}^n -(2(y_i - \mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i)$$



# Interpreting Gradients: 1 Sample

$$L(\mathbf{w}) = \lambda \|\mathbf{w}\|_2^2 + (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

**Recall:**  $\mathbf{w} = \mathbf{w} + -\nabla_{\mathbf{w}}L(\mathbf{w})$  #update  $\mathbf{w}$

$$\nabla_{\mathbf{w}}L(\mathbf{w}) = 2\lambda\mathbf{w} + -(2(y - \mathbf{w}^T \mathbf{x})\mathbf{x})$$

Push  $\mathbf{w}$  towards 0

$$-\nabla_{\mathbf{w}}L(\mathbf{w}) = \underbrace{-2\lambda\mathbf{w}} + \underbrace{(2(y - \mathbf{w}^T \mathbf{x})\mathbf{x})}_{\alpha}$$

If  $y > \mathbf{w}^T \mathbf{x}$  (too low): then  $\mathbf{w} = \mathbf{w} + \alpha \mathbf{x}$  for some  $\alpha$

**Before:**  $\mathbf{w}^T \mathbf{x}$

**After:**  $(\mathbf{w} + \alpha \mathbf{x})^T \mathbf{x} = \mathbf{w}^T \mathbf{x} + \alpha \mathbf{x}^T \mathbf{x}$

# Computing Gradients

- **Numeric gradient:** approximate, slow, easy to write
- **Analytic gradient:** exact, fast, error-prone

In practice: Always use analytic gradient, but check implementation with numerical gradient. This is called a **gradient check**.

```
torch.autograd.gradcheck(func, inputs, eps=1e-06, atol=1e-05, rtol=0.001, raise_exception=True, check_sparse_nnz=False, nondet_tol=0.0)
```

[SOURCE] 

Check gradients computed via small finite differences against analytical gradients w.r.t. tensors in `inputs` that are of floating point type and with `requires_grad=True`.

The check between numerical and analytical gradients uses `allclose()`.

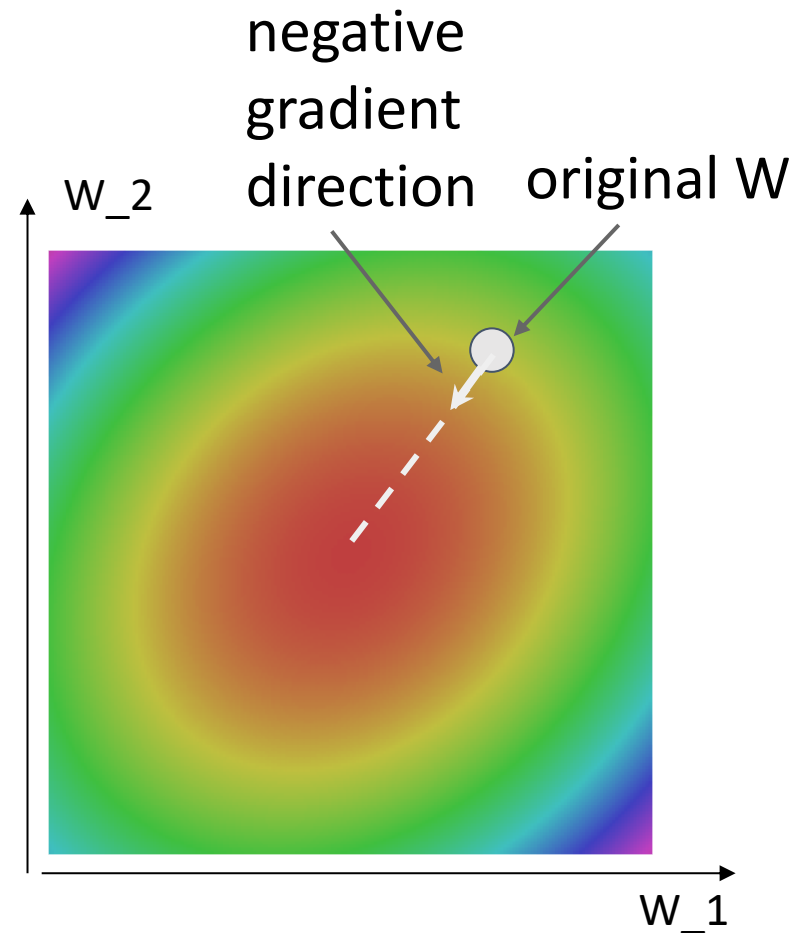
# Gradient Descent

Iteratively step in the direction of the negative gradient  
(direction of local steepest descent)

```
# Vanilla gradient descent
w = initialize_weights()
for t in range(num_steps):
    dw = compute_gradient(loss_fn, data, w)
    w -= learning_rate * dw
```

## Hyperparameters:

- Weight initialization method
- Number of steps
- Learning rate



# Batch Gradient Descent

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

**Problem:** Full sum is expensive when  $N$  is large!

**Solution:** Approximate sum using a minibatch of examples, e.g. 32

# Stochastic Gradient Descent (SGD)

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

**Problem:** Full sum is expensive when  $N$  is large!

**Solution:** Approximate sum using a minibatch of examples, e.g. 32

## Hyperparameters:

- Weight initialization
- Number of steps
- Learning rate
- Batch size
- Data sampling

```
# Stochastic gradient descent
w = initialize_weights()
for t in range(num_steps):
    minibatch = sample_data(data, batch_size)
    dw = compute_gradient(loss_fn, minibatch, w)
    w -= learning_rate * dw
```

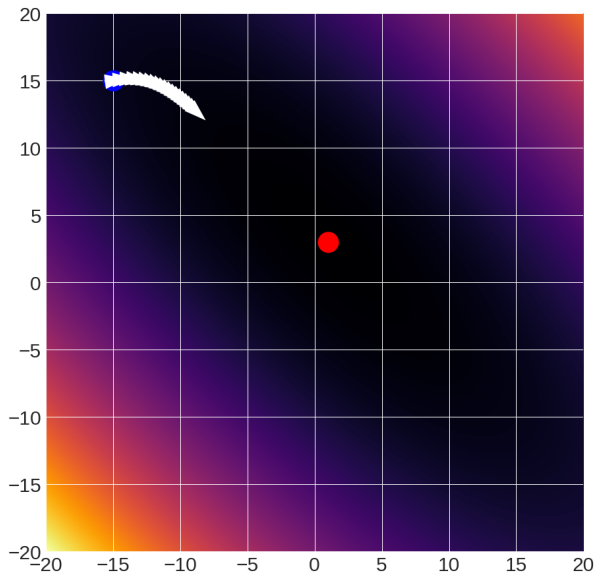
**Note:** Some people say “stochastic gradient descent” is batch size 1, and “minibatch gradient descent” for other batch sizes. I think this distinction is confusing, and use “stochastic gradient descent” for any minibatch size

# Gradient Descent: Learning Rate

Step size (also called **learning rate / lr**)  
*critical parameter*

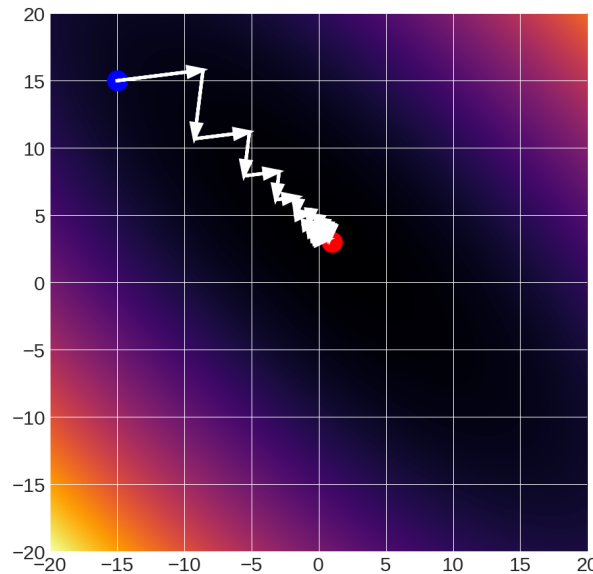
$1 \times 10^{-2}$

falls short



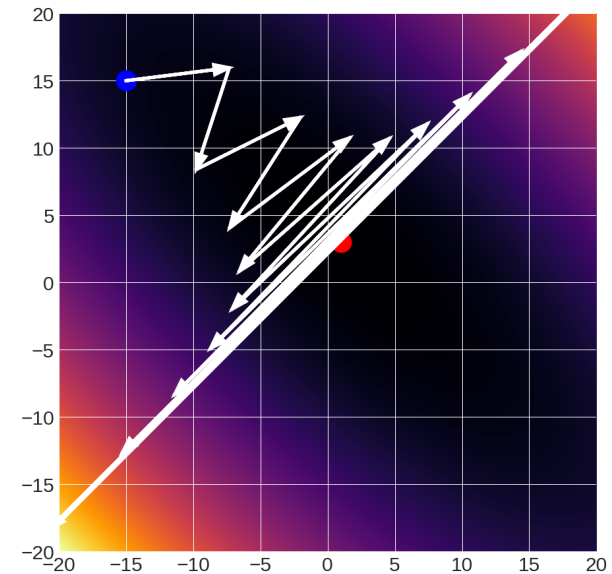
$10 \times 10^{-2}$

converges



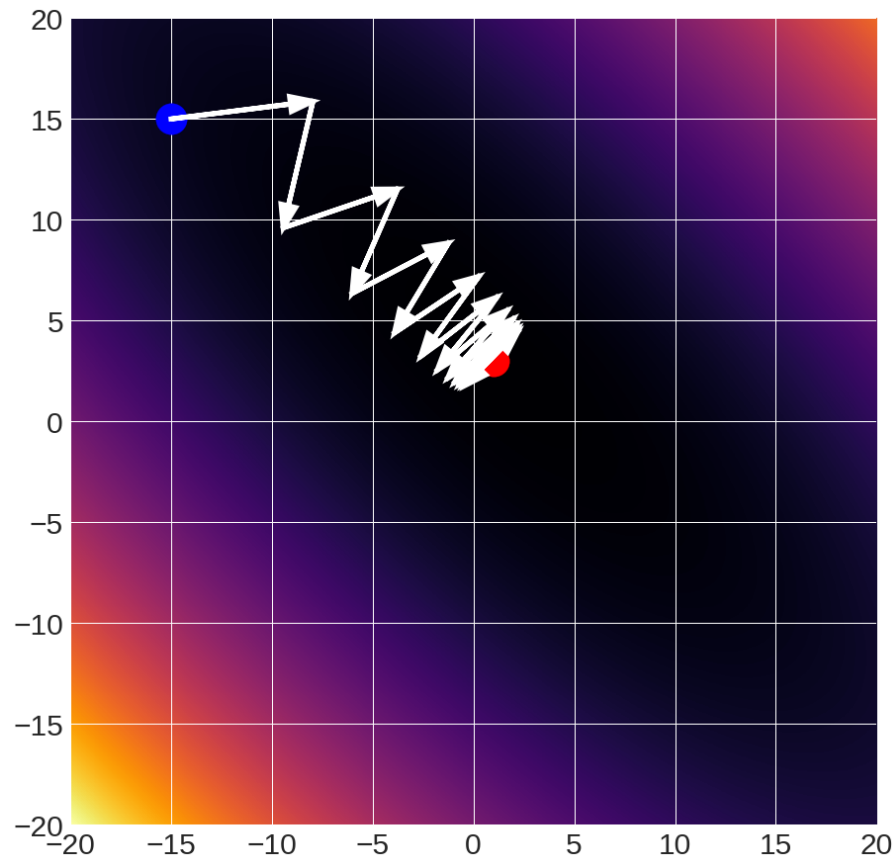
$12 \times 10^{-2}$

diverges



# Gradient Descent: Learning Rate

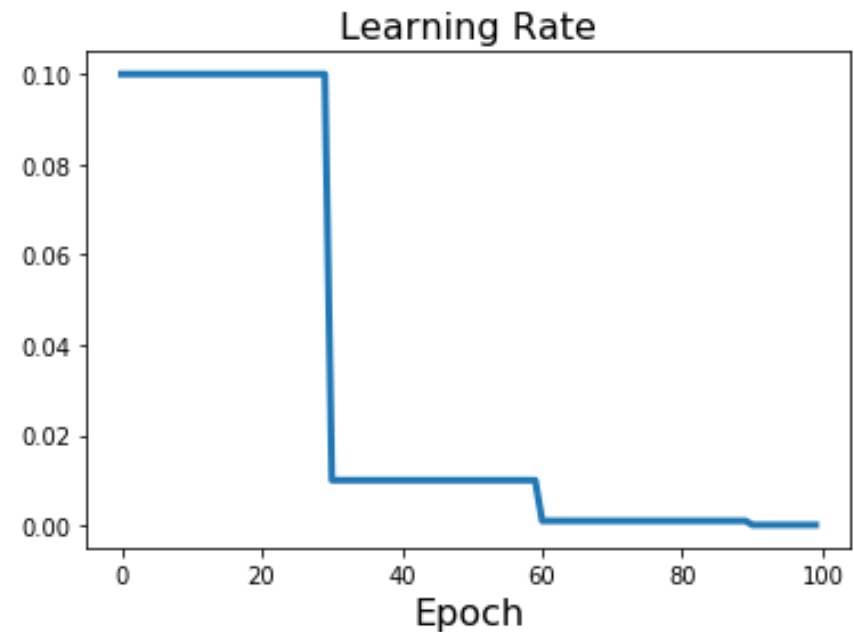
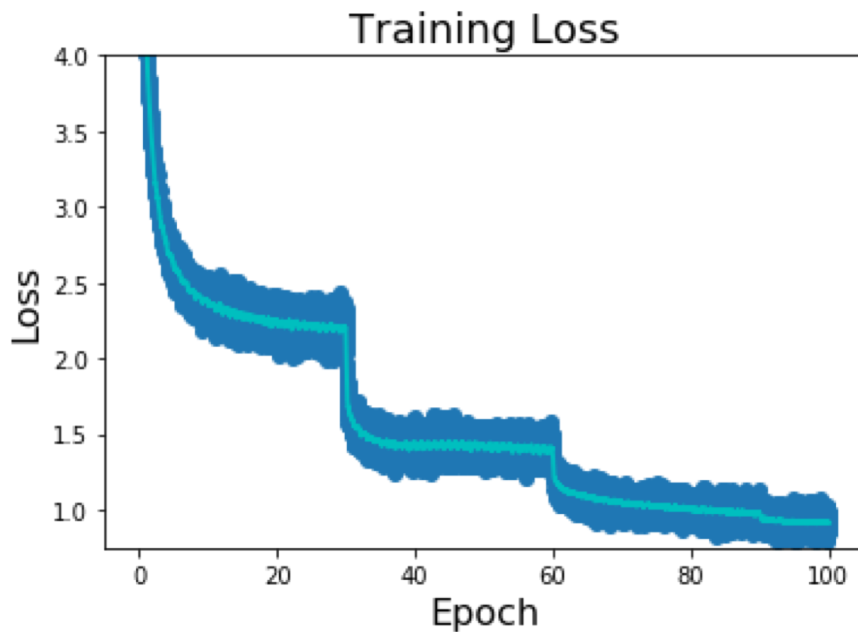
$11 \times 10^{-2}$  : oscillates  
(Raw gradients)



# Learning Rate Decay

**Idea:** Start with high learning rate, reduce it over time.

**Step Decay:** Reduce by some factor at fixed iterations

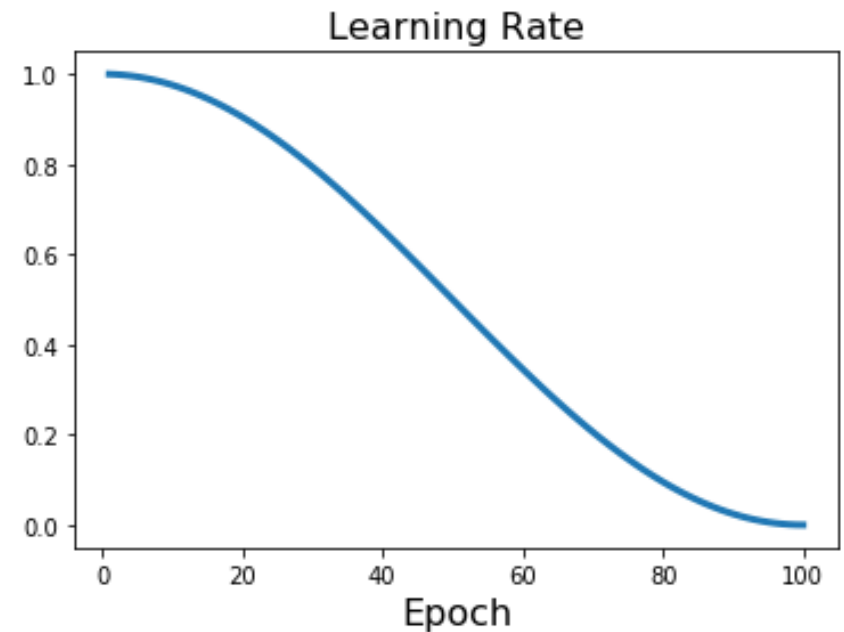
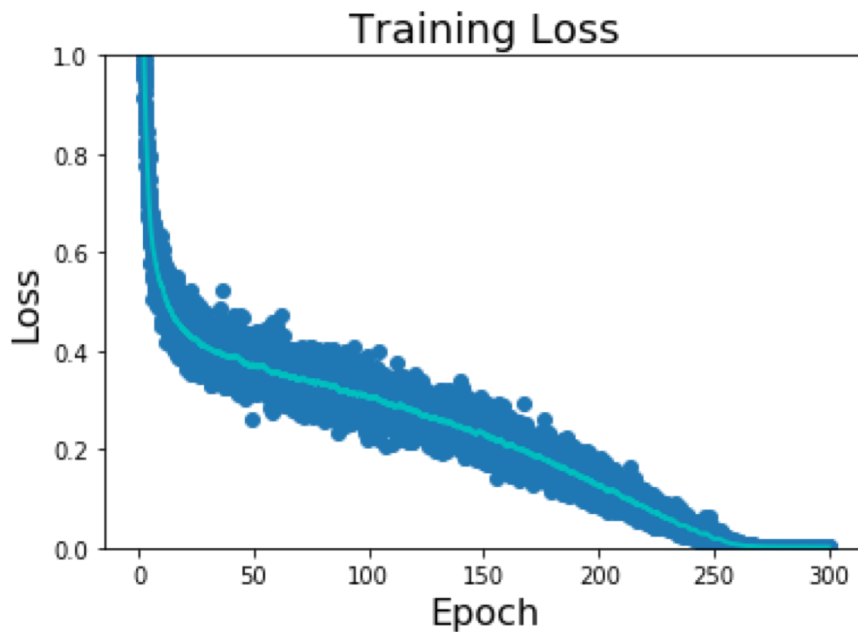




# Learning Rate Decay

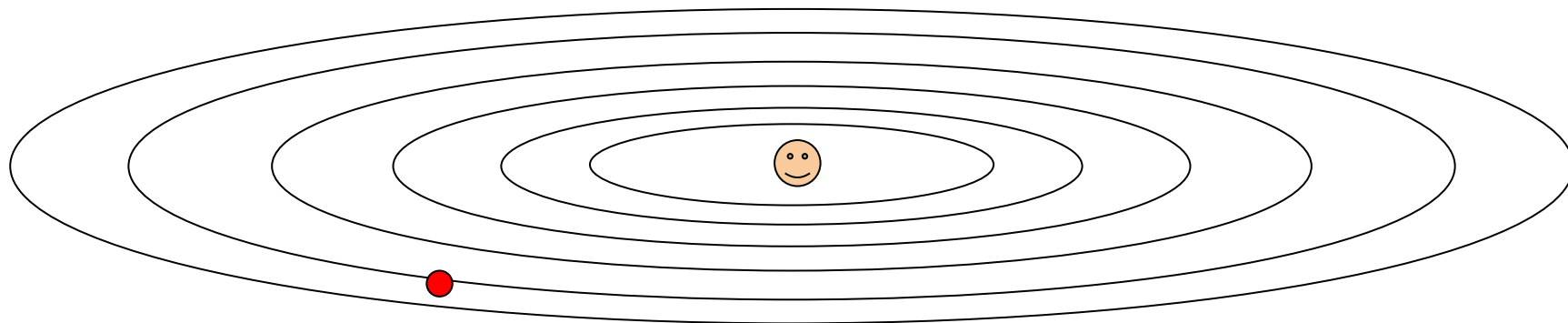
**Idea:** Start with high learning rate, reduce it over time.

**Cosine Decay:** 
$$\alpha_t = \frac{1}{2} \alpha_0 \left( 1 + \cos \left( \frac{t\pi}{T} \right) \right)$$



# Problems with SGD

What if loss changes quickly in one direction and slowly in another?

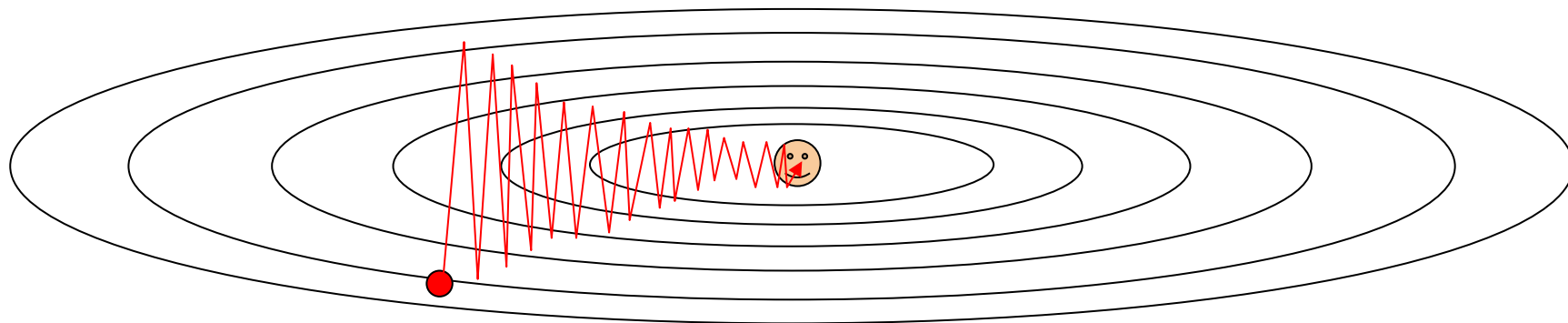


Loss function has high **condition number**: ratio of largest to smallest singular value of the Hessian matrix is large

# Problems with SGD

What if loss changes quickly in one direction and slowly in another?

Slow progress along shallow dimension, jitter along steep direction

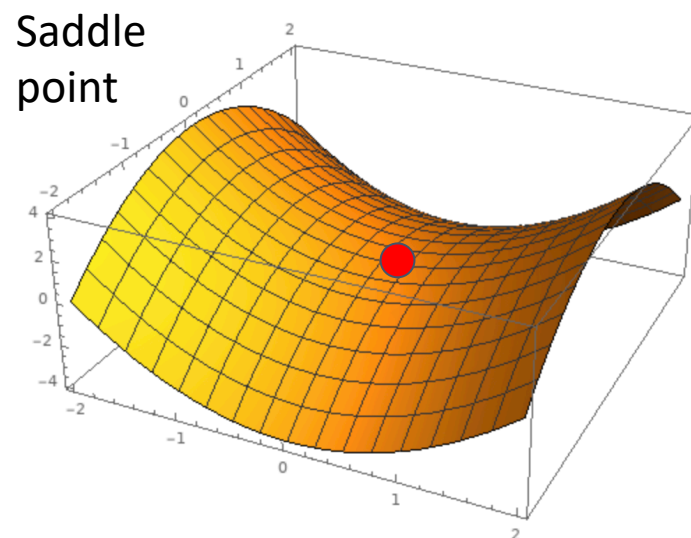


Loss function has high **condition number**: ratio of largest to smallest singular value of the Hessian matrix is large

# Problems with SGD

What if the loss function has a **local minimum** or **saddle point**?

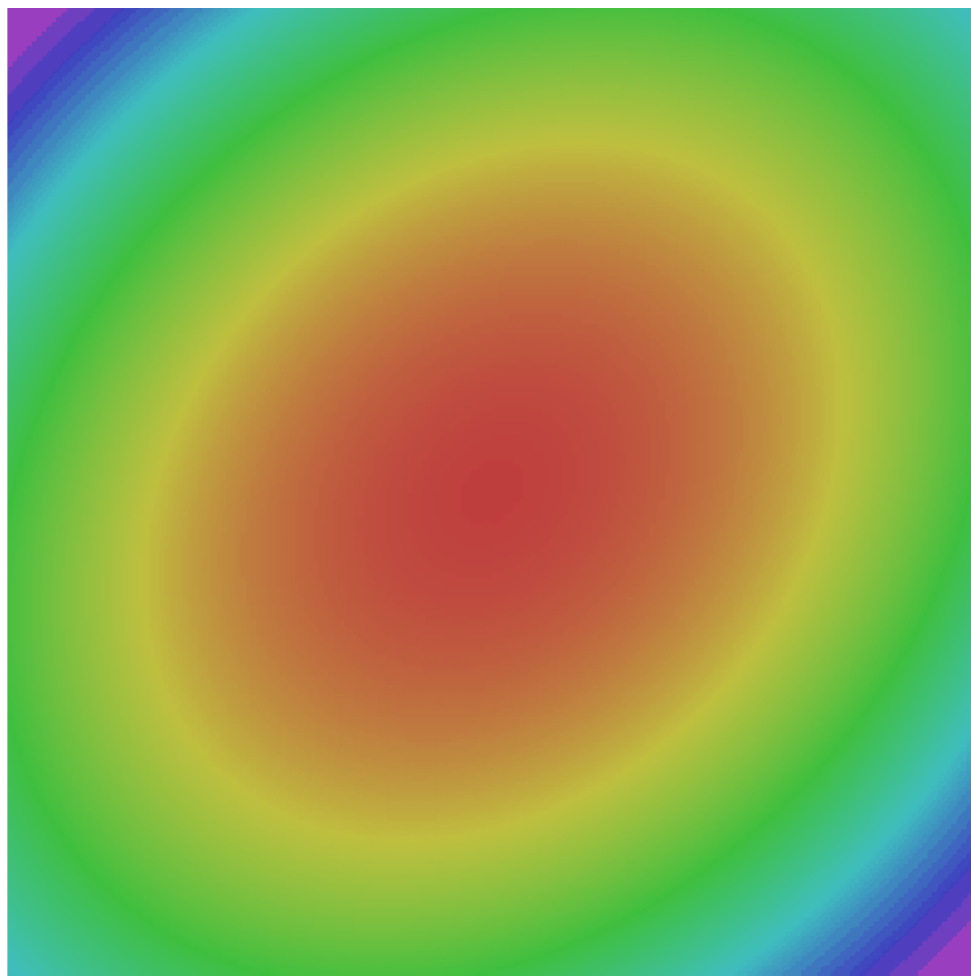
Gradient is zero,  
SGD gets stuck



# Problems with SGD

Our gradients come from minibatches so they can be noisy!

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$



# SGD

## SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

```
for t in range(num_steps):  
    dw = compute_gradient(w)  
    w -= learning_rate * dw
```

# SGD + Momentum

## SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

```
for t in range(num_steps):  
    dw = compute_gradient(w)  
    w -= learning_rate * dw
```

## SGD + Momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$
$$x_{t+1} = x_t - \alpha v_{t+1}$$

```
v = 0  
for t in range(num_steps):  
    dw = compute_gradient(w)  
    v = rho * v + dw  
    w -= learning_rate * v
```

- Build up “velocity” as a running mean of gradients
- Rho gives “friction”; typically  $\rho = 0.9$  or  $0.99$

# SGD + Momentum

## SGD + Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t)$$
$$x_{t+1} = x_t + v_{t+1}$$

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v - learning_rate * dw
    w += v
```

## SGD + Momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$
$$x_{t+1} = x_t - \alpha v_{t+1}$$

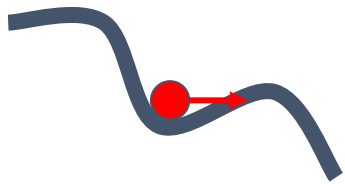
```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v + dw
    w -= learning_rate * v
```

You may see SGD+Momentum formulated different ways, but they are equivalent - give same sequence of  $x$

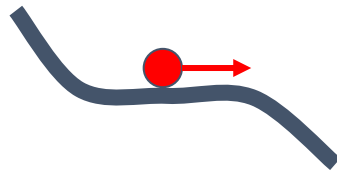


# SGD + Momentum

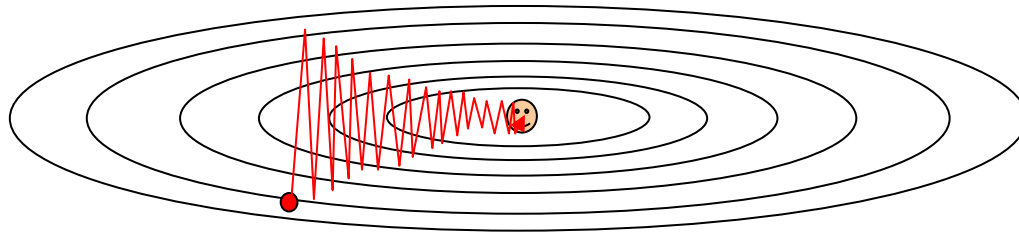
Local Minima



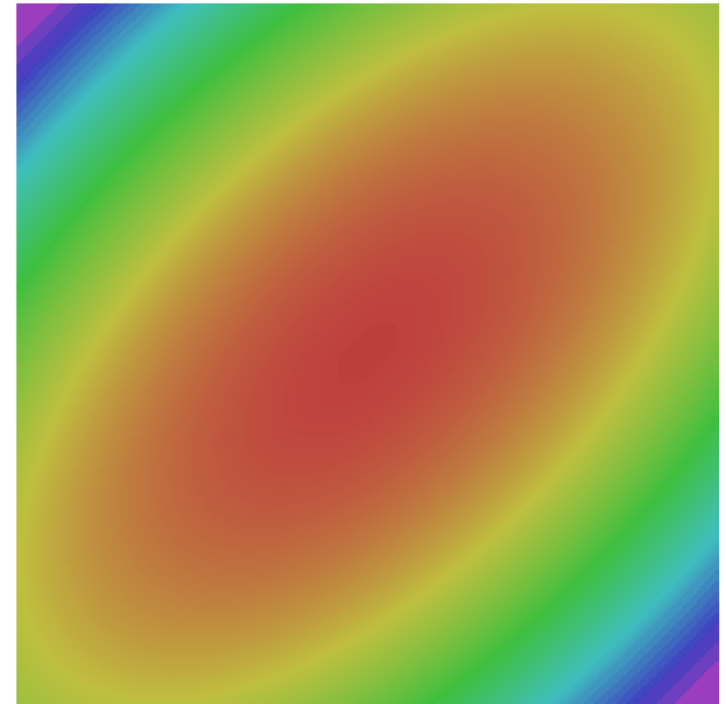
Saddle points



Poor Conditioning



Gradient Noise



— SGD — SGD+Momentum

# Other Update Rules: Adam

```
moment1 = 0
moment2 = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    moment1_unbias = moment1 / (1 - beta1 ** t)
    moment2_unbias = moment2 / (1 - beta2 ** t)
    w -= learning_rate * moment1_unbias / (moment2_unbias.sqrt() + 1e-7)
```

Adam with  $\beta_1 = 0.9$ ,  
 $\beta_2 = 0.999$ , and  $\text{learning\_rate} = 1e-3, 5e-4, 1e-4$   
is a great starting point for many models!

# Adam: Very Common in Practice!

for input to the CNN; each colored pixel in the image yields a 7D one-hot vector. Following common practice, the network is trained end-to-end using stochastic gradient descent with the Adam optimizer [22]. We anneal the learning rate to 0 using a half cosine schedule without restarts [28].

Bakhtin, van der Maaten, Johnson, Gustafson, and Girshick, NeurIPS 2019

We train all models using Adam [23] with learning rate  $10^{-4}$  and batch size 32 for 1 million iterations; training takes about 3 days on a single Tesla P100. For each mini-batch we first update  $f$ , then update  $D_{img}$  and  $D_{obj}$ .

Johnson, Gupta, and Fei-Fei, CVPR 2018

ganized into three residual blocks. We train for 25 epochs using Adam [27] with learning rate  $10^{-4}$  and 32 images per batch on 8 Tesla V100 GPUs. We set the `cubify` thresh-

Gkioxari, Malik, and Johnson, ICCV 2019

sampled with each bit drawn uniformly at random. For gradient descent, we use Adam [29] with a learning rate of  $10^{-3}$  and default hyperparameters. All models are trained with batch size 12. Models are trained for 200 epochs, or 400 epochs if being trained on multiple noise layers.

Zhu, Kaplan, Johnson, and Fei-Fei, ECCV 2018

16 dimensional vectors. We iteratively train the Generator and Discriminator with a batch size of 64 for 200 epochs using Adam [22] with an initial learning rate of 0.001.

Gupta, Johnson, et al, CVPR 2018

Adam with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ , and `learning_rate = 1e-3, 5e-4, 1e-4` is a great starting point for many models!

# Optimization in Practice

- **Conventional wisdom:** minibatch stochastic gradient descent (SGD) + momentum (package implements it for you) + some sensibly changing learning rate
- The above is typically what is meant by “SGD”
- Other update rules exist (Adam very common); sometimes better, sometimes worse than SGD

# Optimizing Everything

$$L(\mathbf{W}) = \lambda \|\mathbf{W}\|_2^2 + \sum_{i=1}^n -\log \left( \frac{\exp((W\mathbf{x})_{y_i})}{\sum_k \exp((W\mathbf{x})_k)} \right)$$

$$L(\mathbf{w}) = \lambda \|\mathbf{w}\|_2^2 + \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

- Optimize  $\mathbf{w}$  on training set with SGD to maximize training accuracy
- Optimize  $\lambda$  with random/grid search to maximize validation accuracy
- Note: Optimizing  $\lambda$  on training sets it to 0

# Overfitting / Underfitting and Model Complexity

# (Over/Under)fitting and Complexity

Let's fit a polynomial: given  $x$ , predict  $y$

$$y = w_0 + w_1x + w_2x^2 + w_3x^3 + \dots + w_Fx^F$$

Note: can do non-linear regression with copies of  $x$

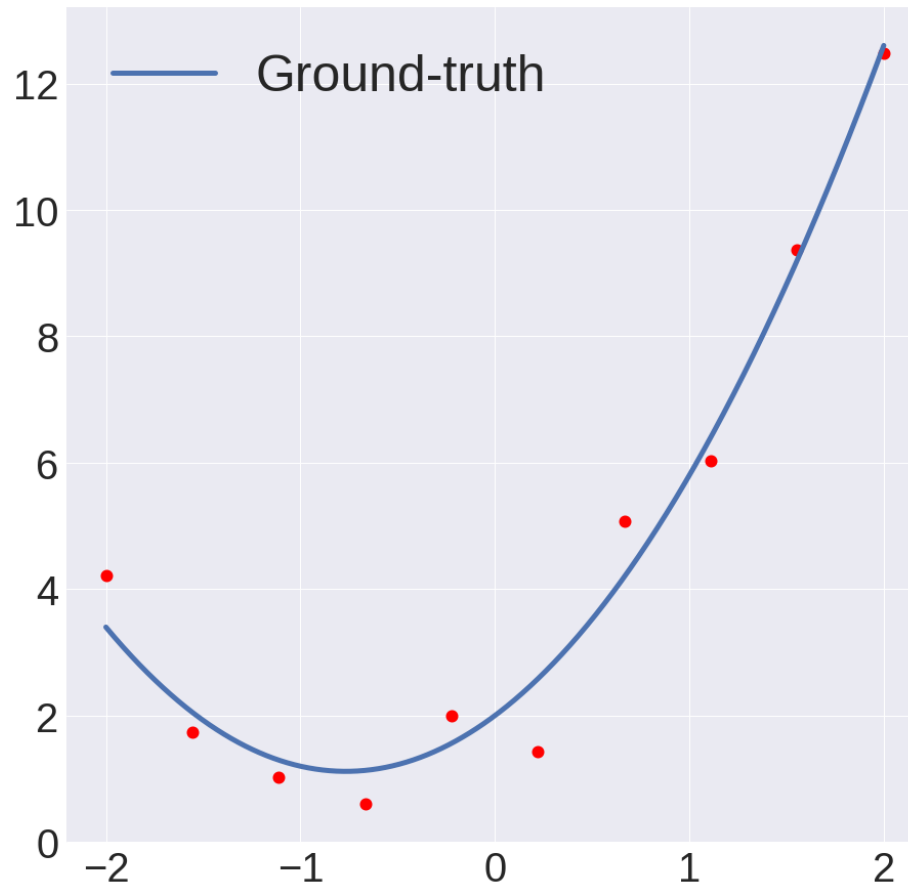
$$\begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} x_1^F & \cdots & x_1^2 & x_1 & 1 \\ \vdots & \ddots & \vdots & \vdots & \vdots \\ x_N^F & \cdots & x_N^2 & x_N & 1 \end{bmatrix} \begin{bmatrix} w_F \\ \vdots \\ w_2 \\ w_1 \\ w_0 \end{bmatrix}$$

Matrix of all polynomial degrees ↑

Weights: one per polynomial degree ↑

# (Over/Under)fitting and Complexity

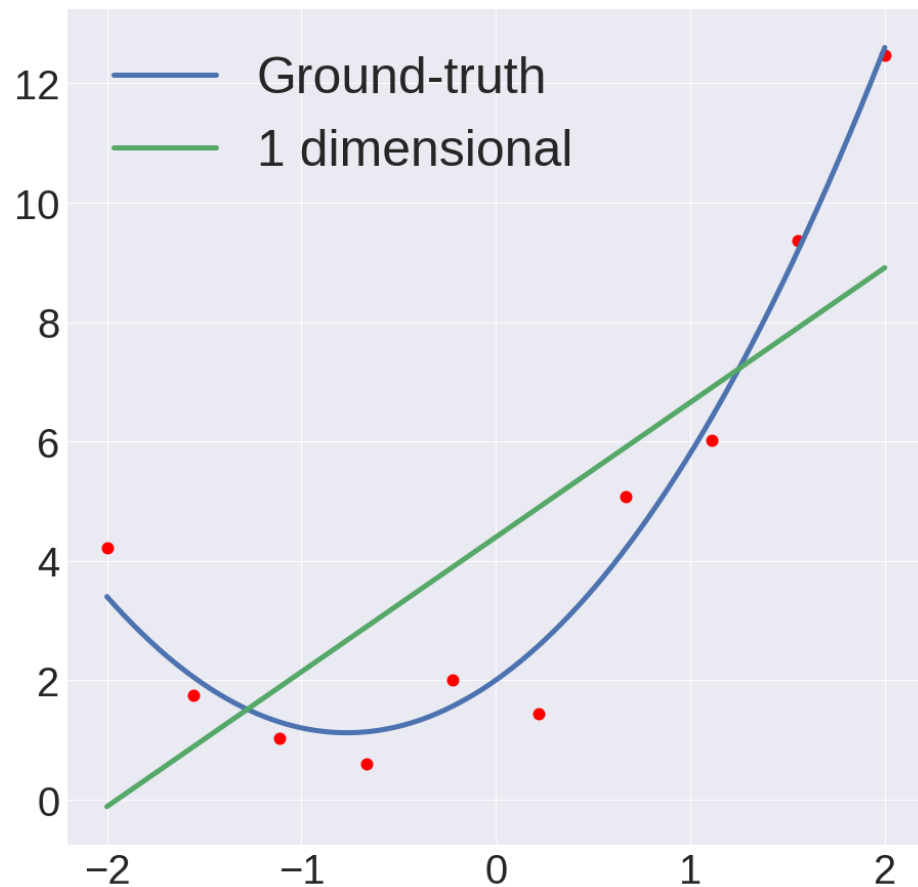
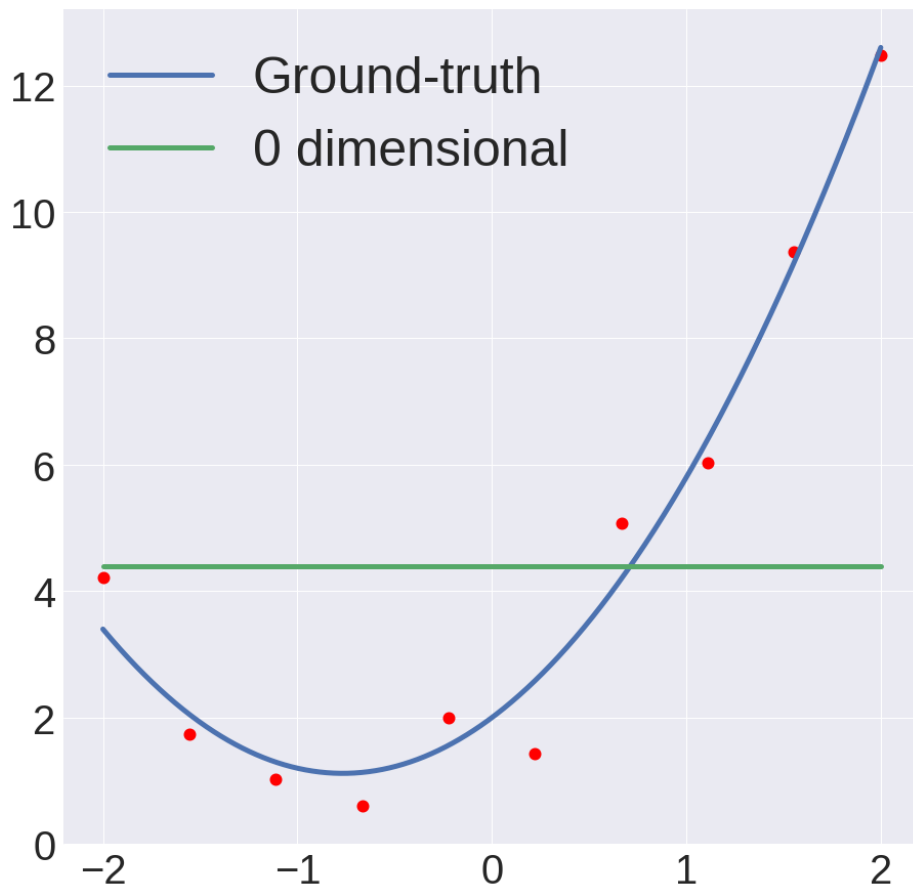
**Ground-Truth:**  $1.5x^2 + 2.3x + 2 + N(0,0.5)$





# Underfitting

**Ground-Truth:  $1.5x^2 + 2.3x + 2 + N(0,0.5)$**

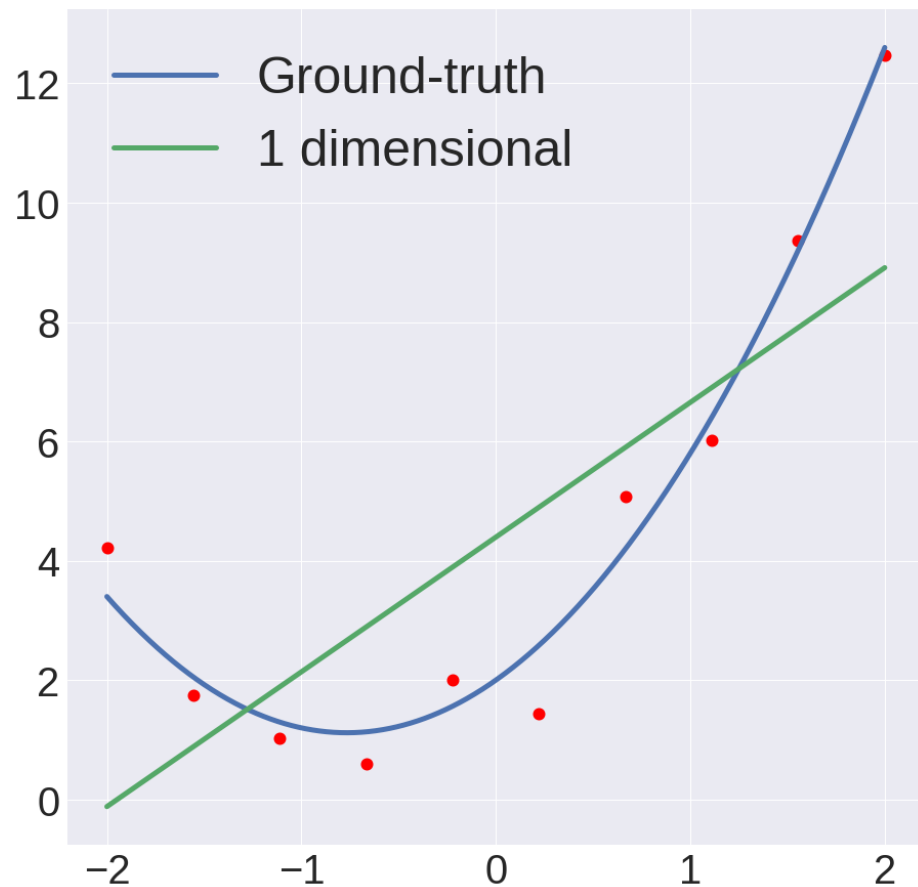


# Underfitting

**Ground-Truth:**  $1.5x^2 + 2.3x + 2 + N(0,0.5)$

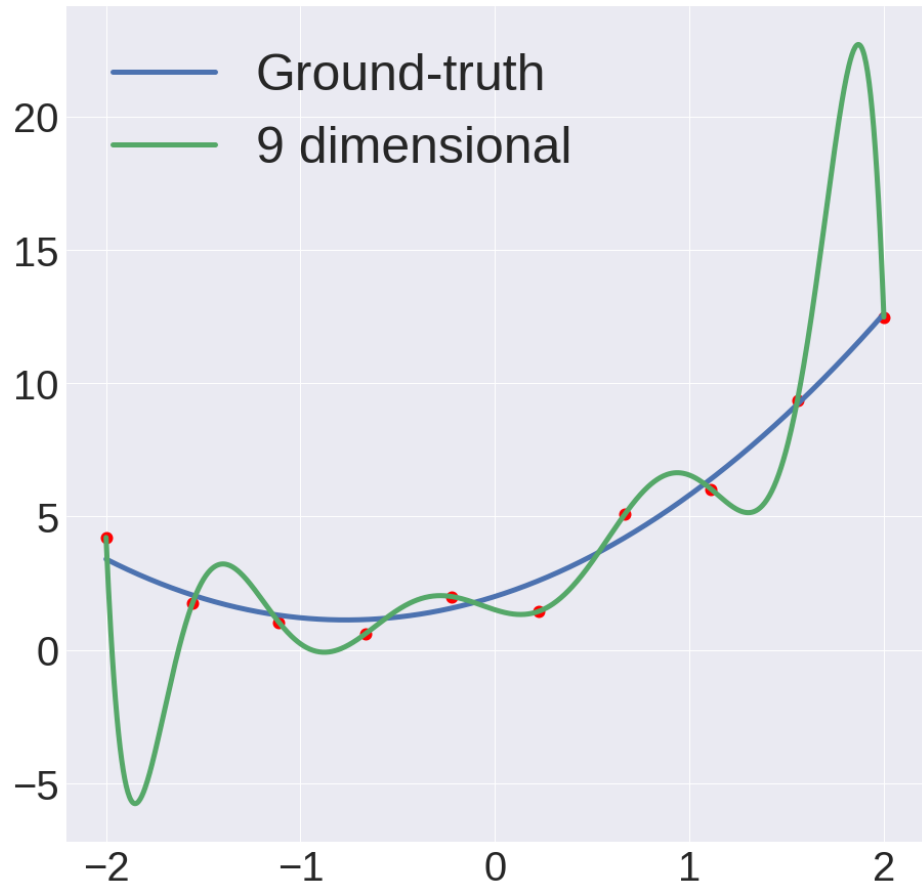
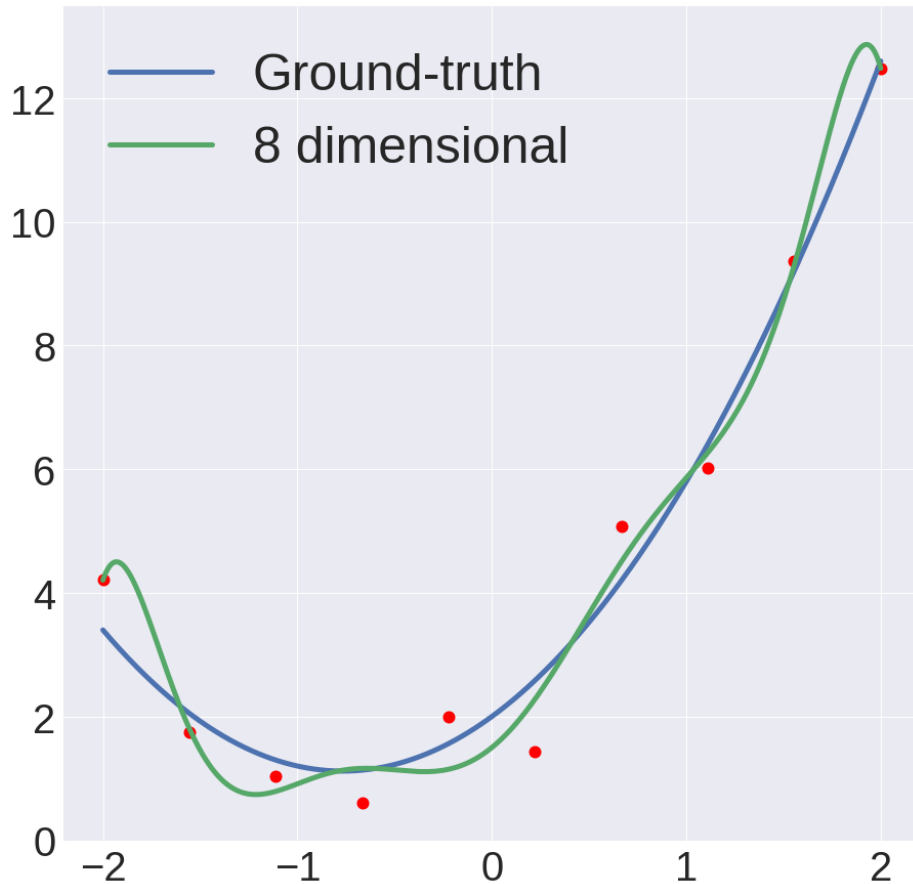
Model isn't "complex" enough to fit the data

*Bias* (statistics): Error intrinsic to the model.



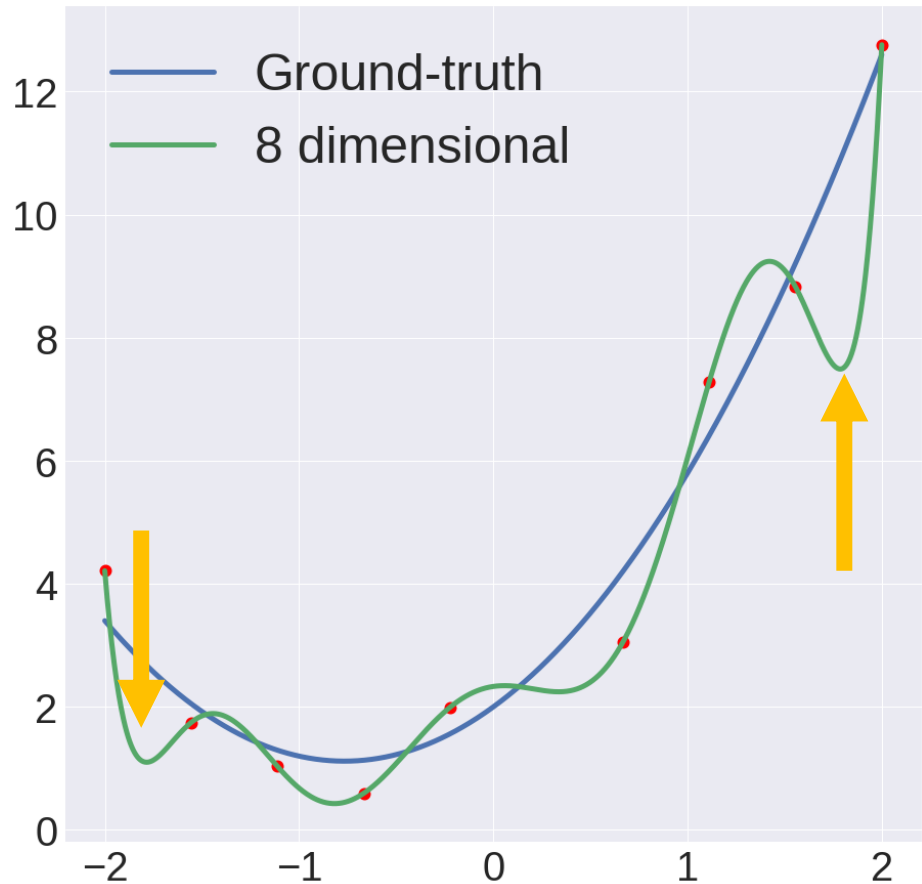
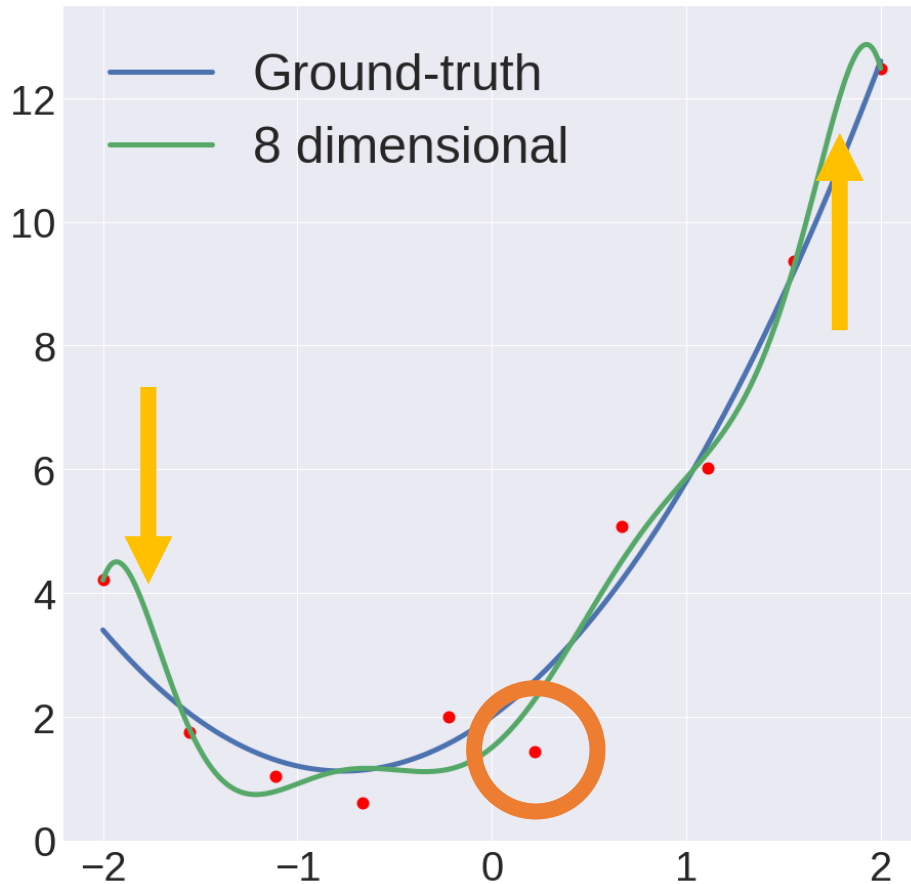
# Overfitting

**Ground-Truth:  $1.5x^2 + 2.3x + 2 + N(0,0.5)$**




# Overfitting

Model has high ***variance***: remove **one point**, and model changes dramatically



# (Continuous) Model Complexity

$$\arg \min_W \lambda \|W\|_2^2 + \sum_{i=1}^n \underbrace{-\log \left( \frac{\exp((Wx)_{y_i})}{\sum_k \exp((Wx)_k)} \right)}_{\text{Pay penalty for negative log-likelihood of correct class}}$$


Regularization: penalty  
for complex model

Pay penalty for negative log-  
likelihood of correct class

Intuitively: big weights = more complex model

Model 1:  $0.01 * x_1 + 1.3 * x_2 + -0.02 * x_3 + -2.1x_4 + 10$

Model 2:  $37.2 * x_1 + 13.4 * x_2 + 5.6 * x_3 + -6.1x_4 + 30$

# Fitting a Model

Again, fitting polynomial, but with regularization

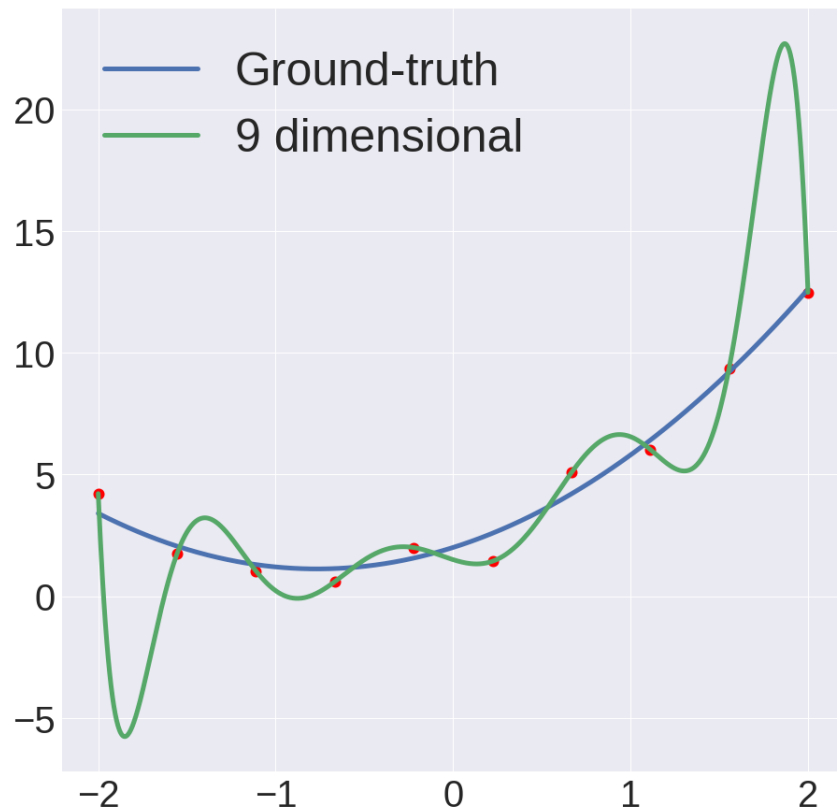
$$\arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\| + \lambda \|\mathbf{w}\|$$

The diagram illustrates the components of the optimization problem. A blue arrow points from the  $\arg \min$  operator to the feature matrix  $\mathbf{X}$ . Another blue arrow points from the  $\mathbf{X}$  in the first term to the matrix representation. Two orange arrows point from the  $\mathbf{w}$  terms in the first and second terms to the weight vector representation.

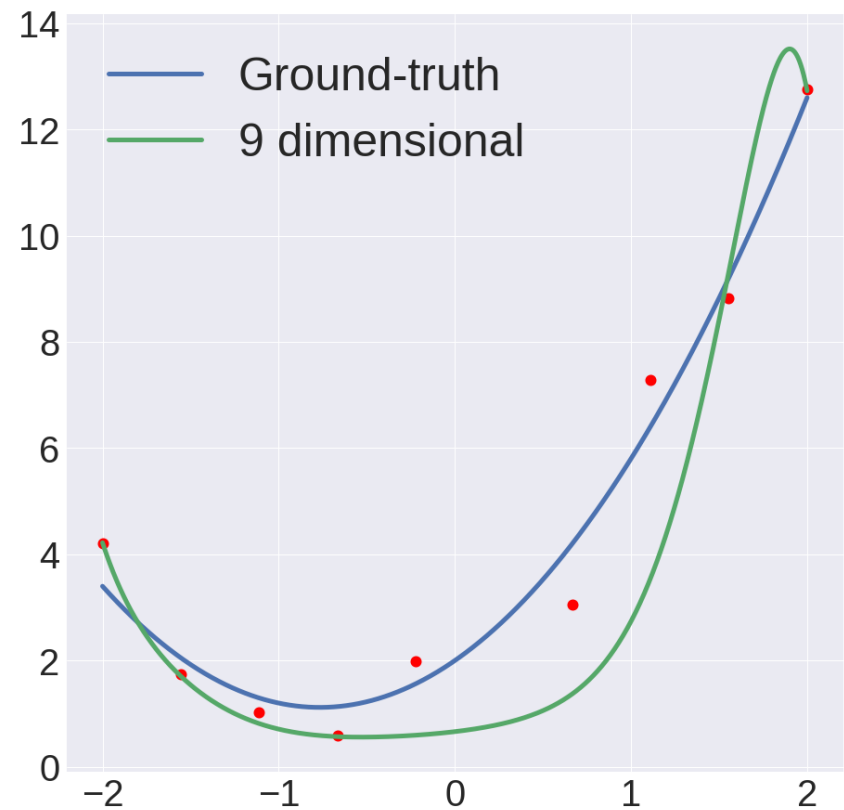
$$\begin{bmatrix} x_1^F & \cdots & x_1^2 & x_1 & 1 \\ \vdots & \ddots & \vdots & \vdots & \vdots \\ x_N^F & \cdots & x_N^2 & x_N & 1 \end{bmatrix} \quad \begin{bmatrix} w_F \\ \vdots \\ w_0 \end{bmatrix}$$

# Adding Regularization

No regularization:  
fits all data points



Regularization:  
can't fit all data points



# Bias / Variance Tradeoff

Error on new data comes from combination of:

- 1. Bias:** model is oversimplified and can't fit the underlying data
- 2. Variance:** you don't have the ability to estimate your model from limited data
- 3. Inherent:** the data is intrinsically difficult

Bias and variance trade-off. Fixing one hurts the other. You can prove theorems about this.



# Underfitting and Overfitting

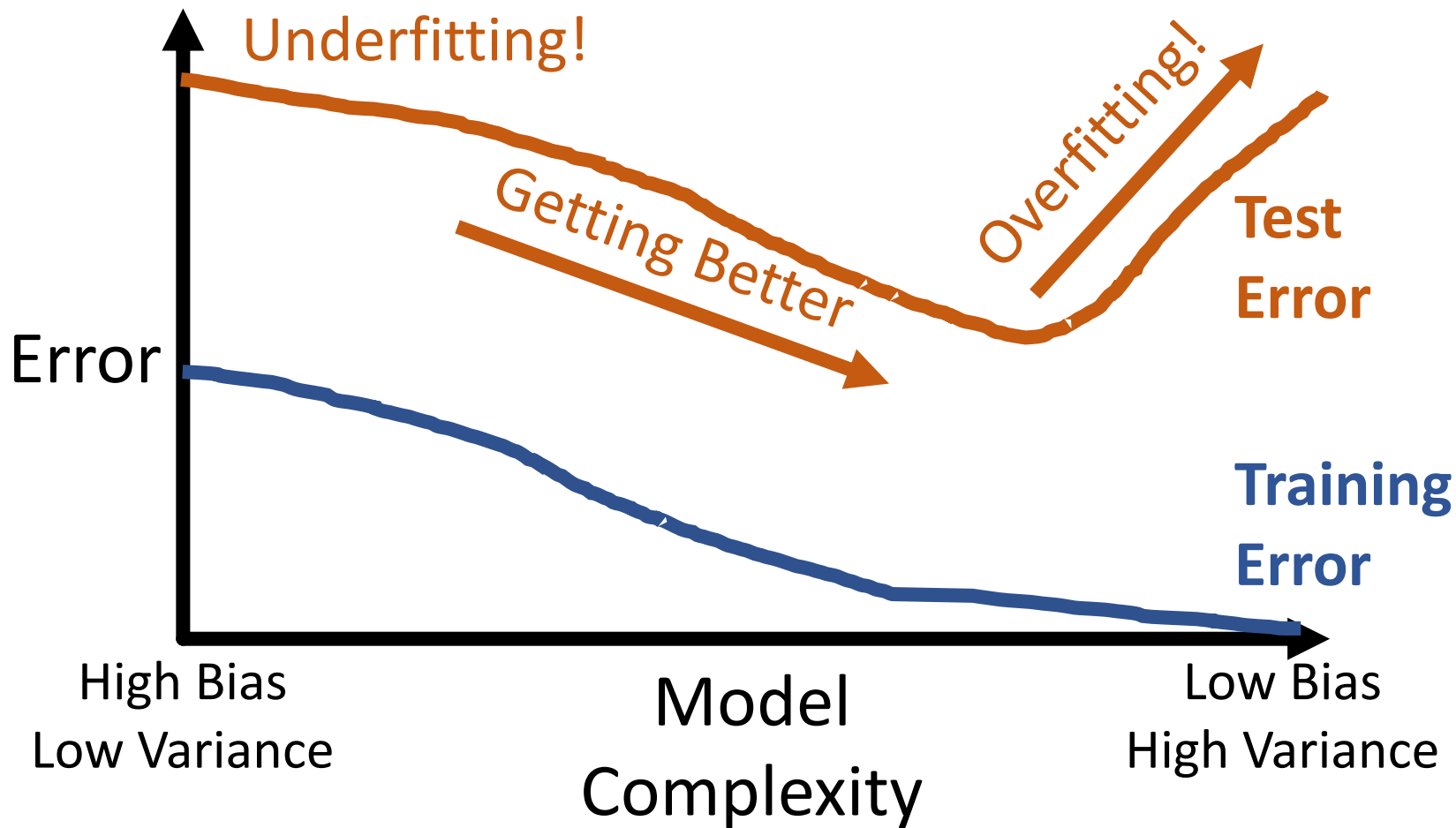


Diagram adapted from: D. Hoiem

# Underfitting and Overfitting

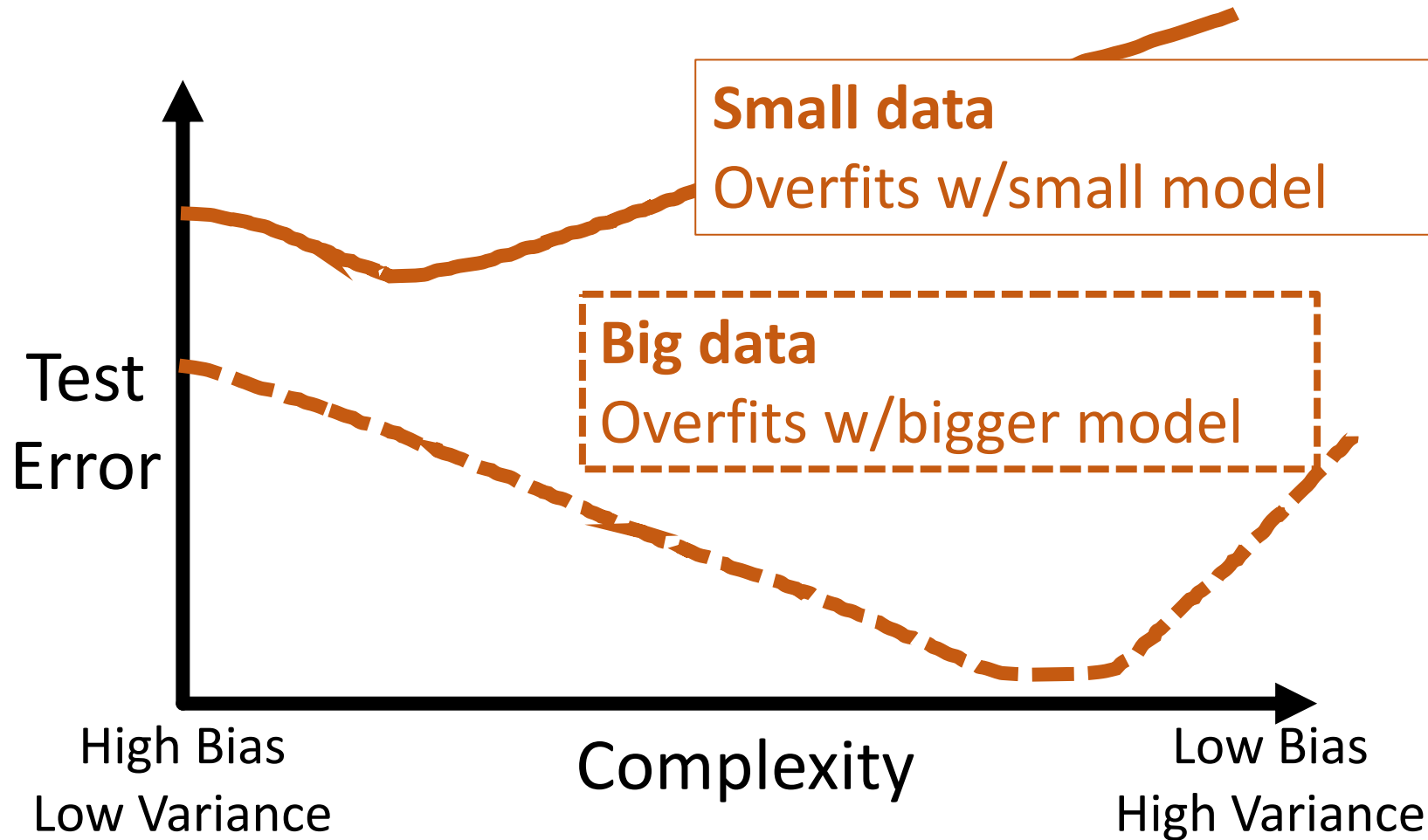
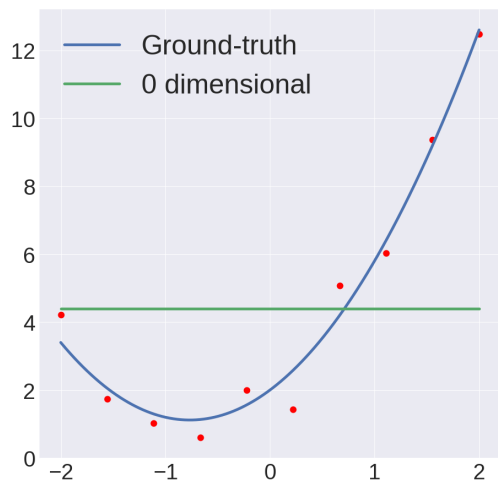


Diagram adapted from: D. Hoiem

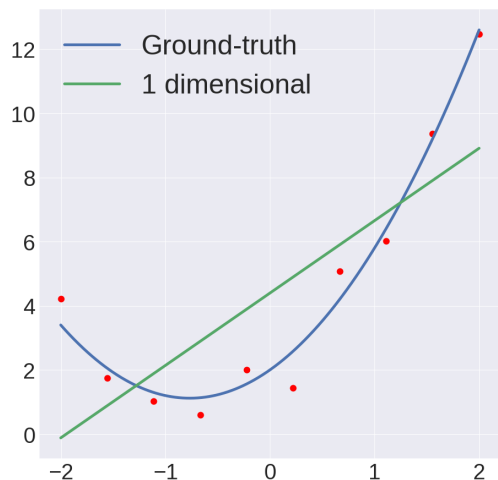
# Underfitting



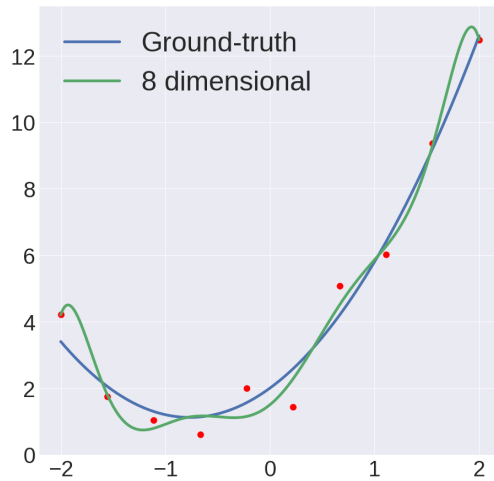
Do poorly on both training and validation data due to bias.

Solution:

1. More features
2. More powerful model
3. Reduce regularization



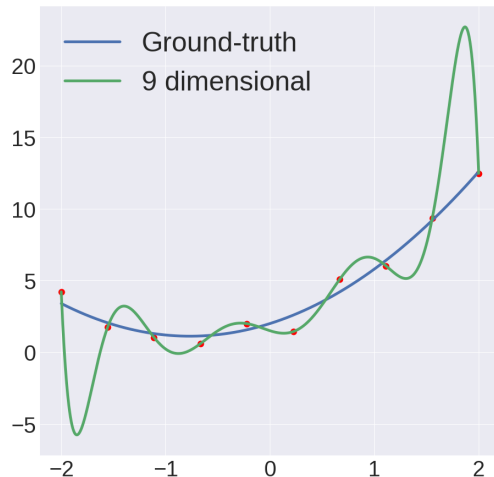
# Overfitting



Do well on training data, but poorly on validation data due to variance

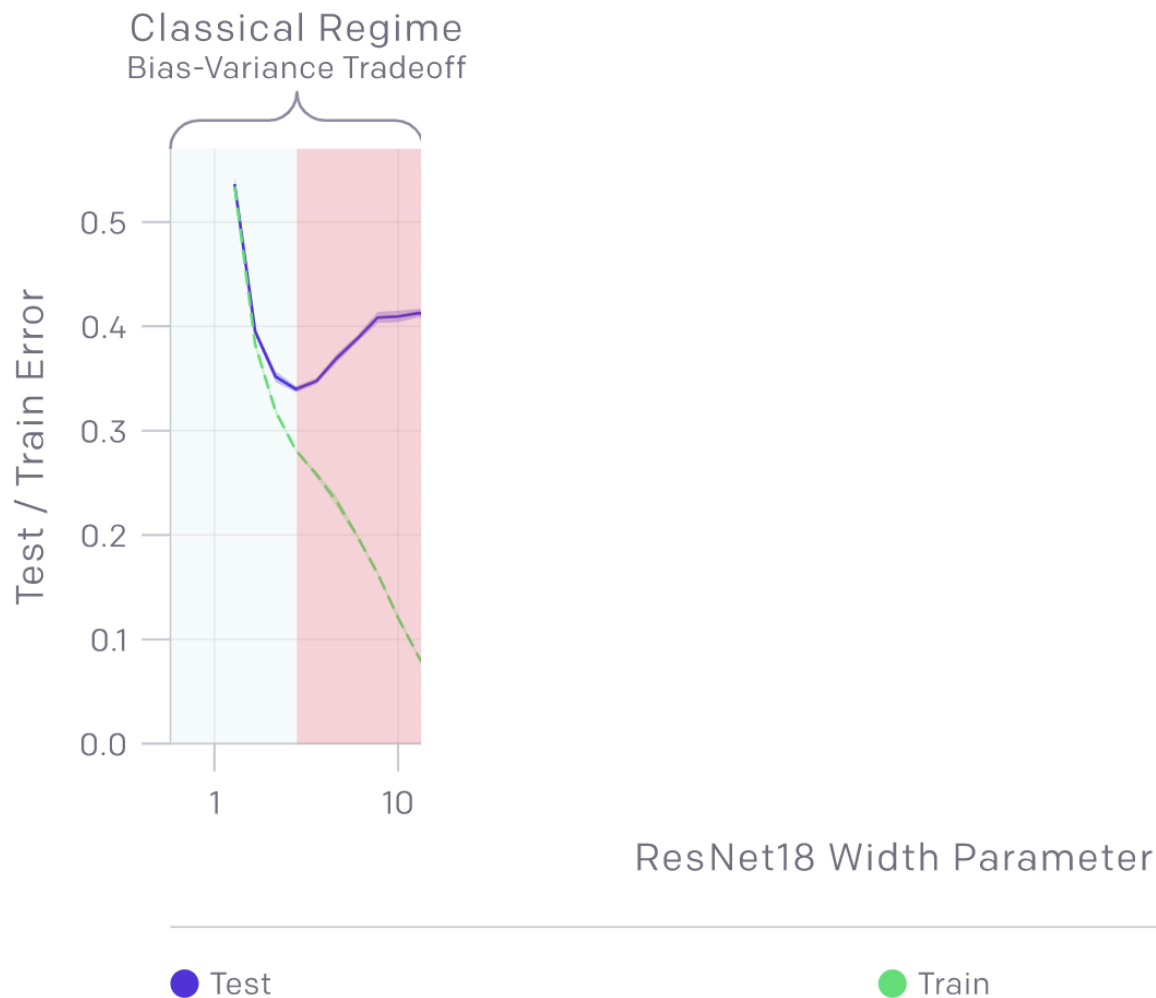
Solution:

1. More data
2. Less powerful model
3. Regularize your model more

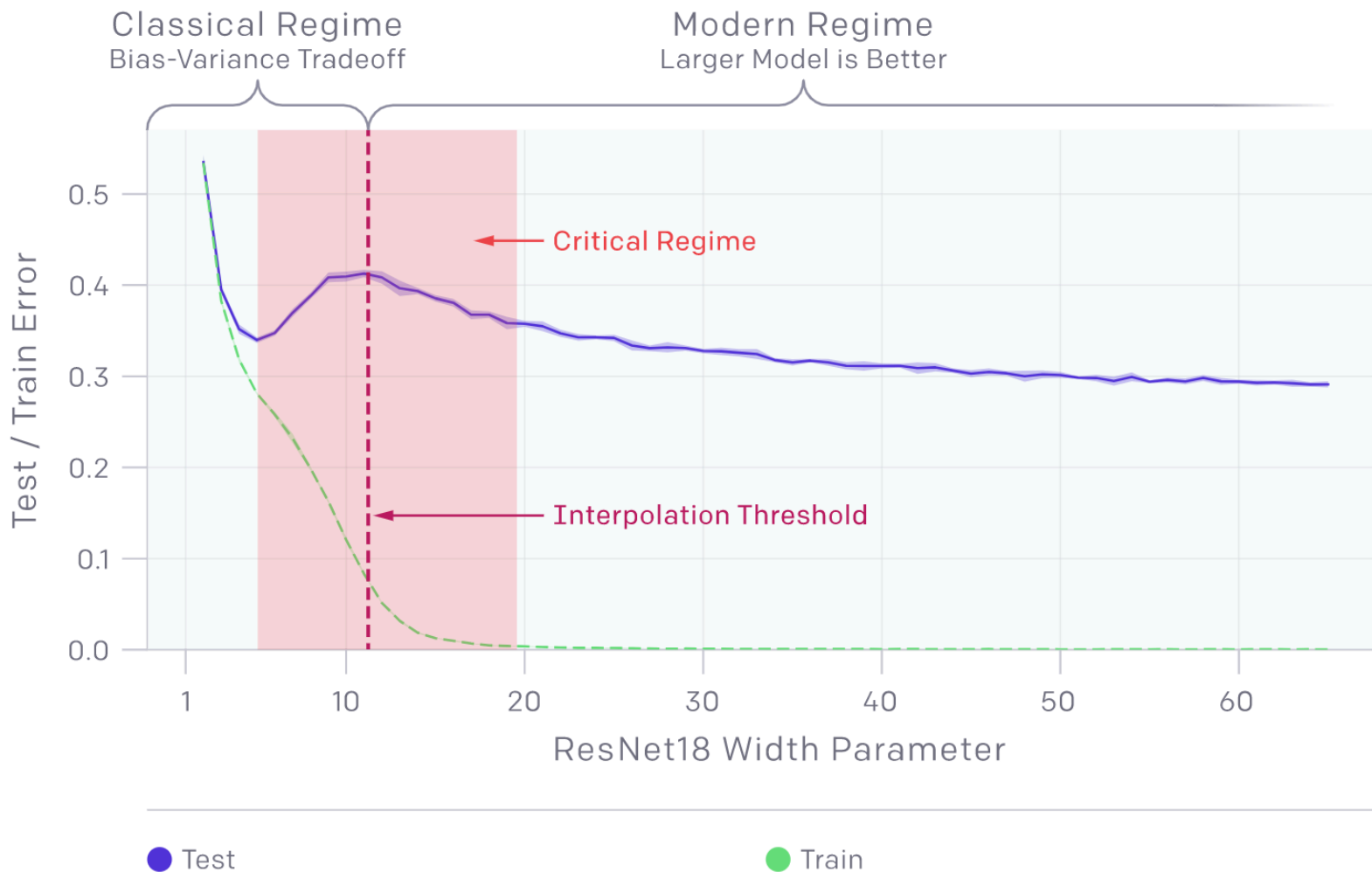


Heuristic: First make sure you *can* overfit, then stop overfitting.

# Double Descent



# Double Descent



Advani and Saxe, "High-dimensional dynamics of generalization error in neural networks", 2017

Geiger et al, "The jamming transition as a paradigm to understand the loss landscape of deep neural networks", 2018

Belkin et al, "Reconciling modern machine learning practice and the bias-variance trade-off", 2018

Nakkiran et al, "Deep Double Descent: Where Bigger Models and More Data Hurt", 2019

# Recap

Next Time:  
Nonlinear Models,  
Neural Networks!