

CMPT 318

Data Analytics

Learning Algorithms II

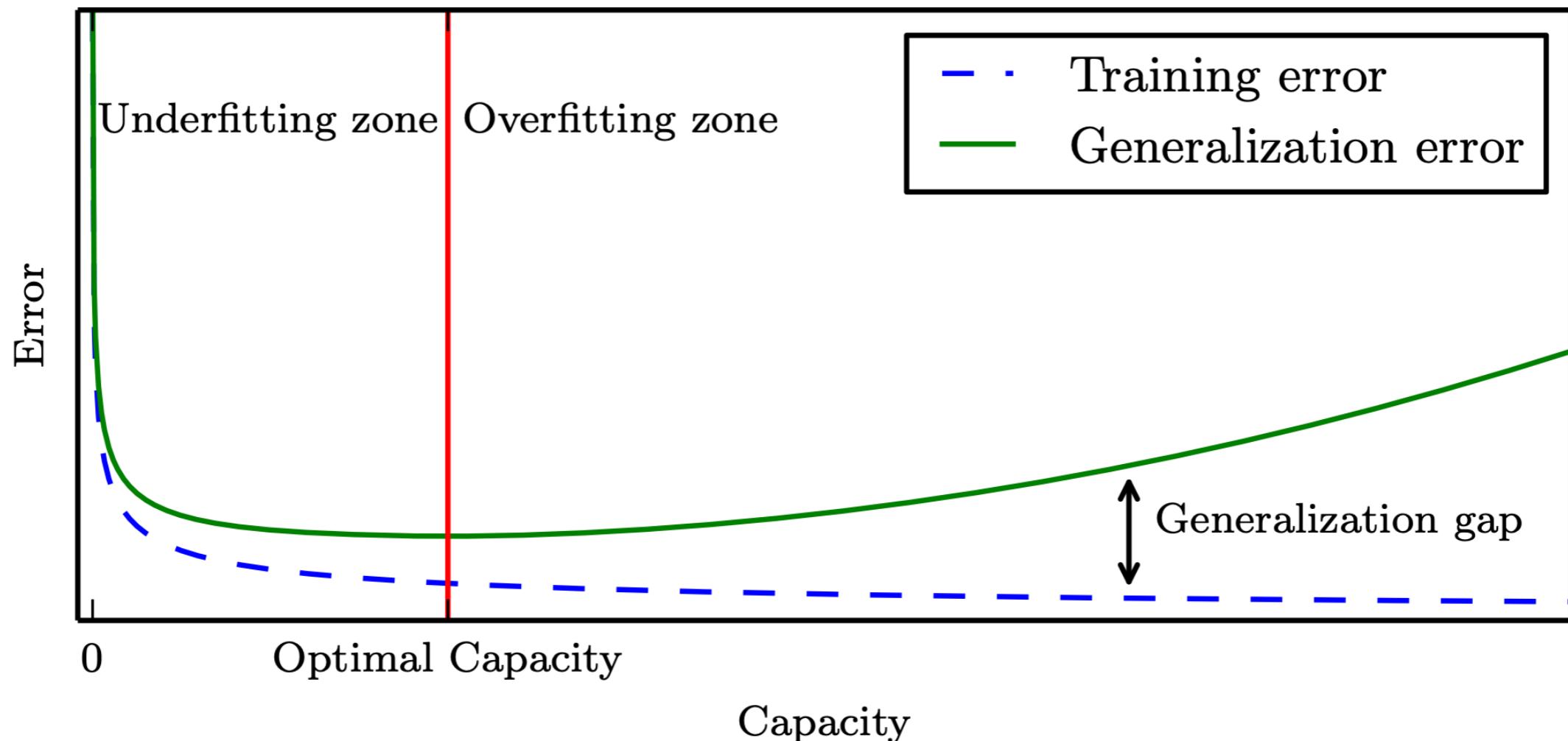
Adapted in part from
Chapter 8 of Data Science from Scratch (Grus)
Chapter 4, 5 of Deep Learning (Goodfellow)
Chapter 1 of Convex Optimization (Boyd, ISBN: 9780521833783)

Upcoming lecture schedule

- Proposals will be graded/returned Thursday
- Friday (Nov. 3) - Project Q&A, individual meetings to discuss any project roadblocks. Attendance optional.
- Monday (Nov. 13) - starting on Neural Networks

Review

Generalization and Capacity



Regularization

Any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.

Regularization with Weight Decay

- We can reduce model capacity by:
 1. Allowing lots of functions, but providing a preference for one solution over another,
e.g. prefer solutions with smaller norm

$$J(\mathbf{w}) = \text{MSE}_{\text{train}} + \lambda \mathbf{w}^T \mathbf{w}, \quad (5.18)$$

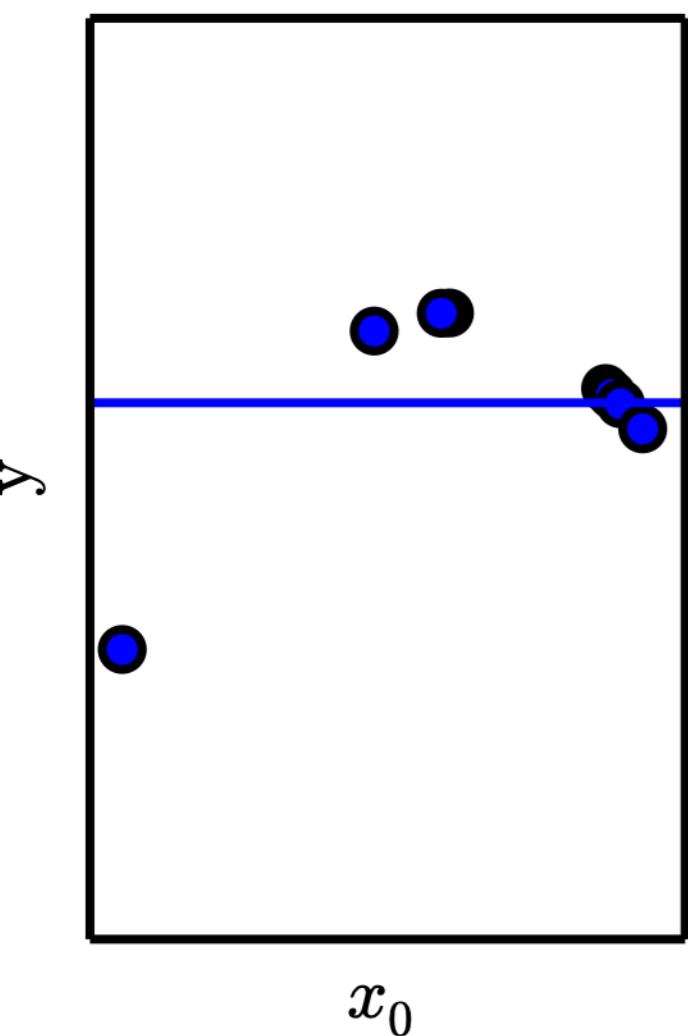
 



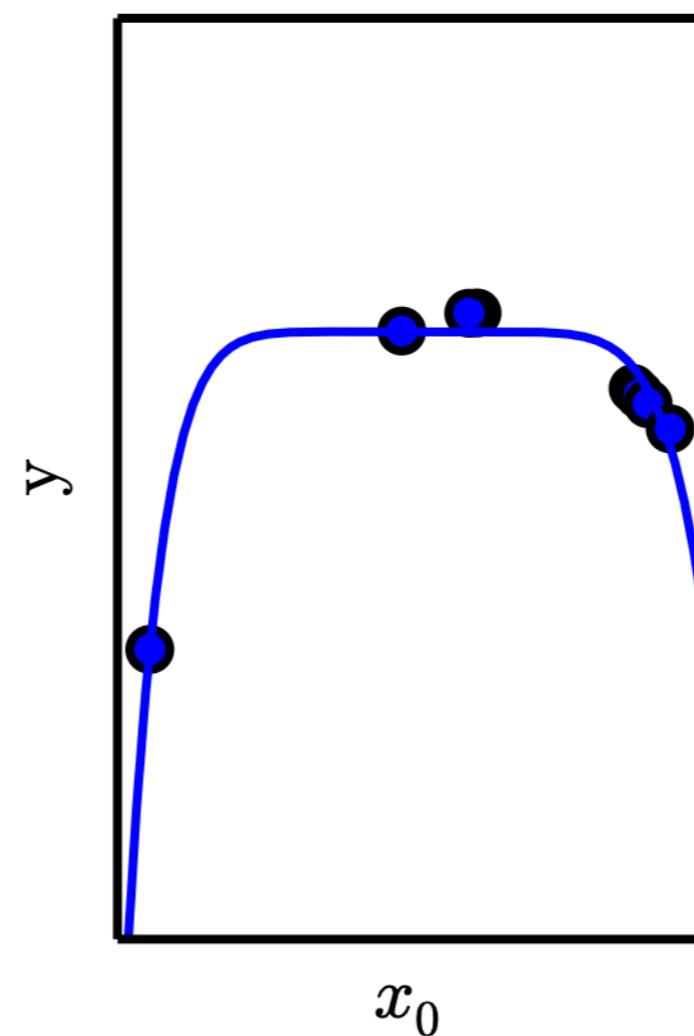
Review

Weight Decay

Underfitting
(Excessive λ)



Appropriate weight decay
(Medium λ)



Overfitting
($\lambda \rightarrow 0$)

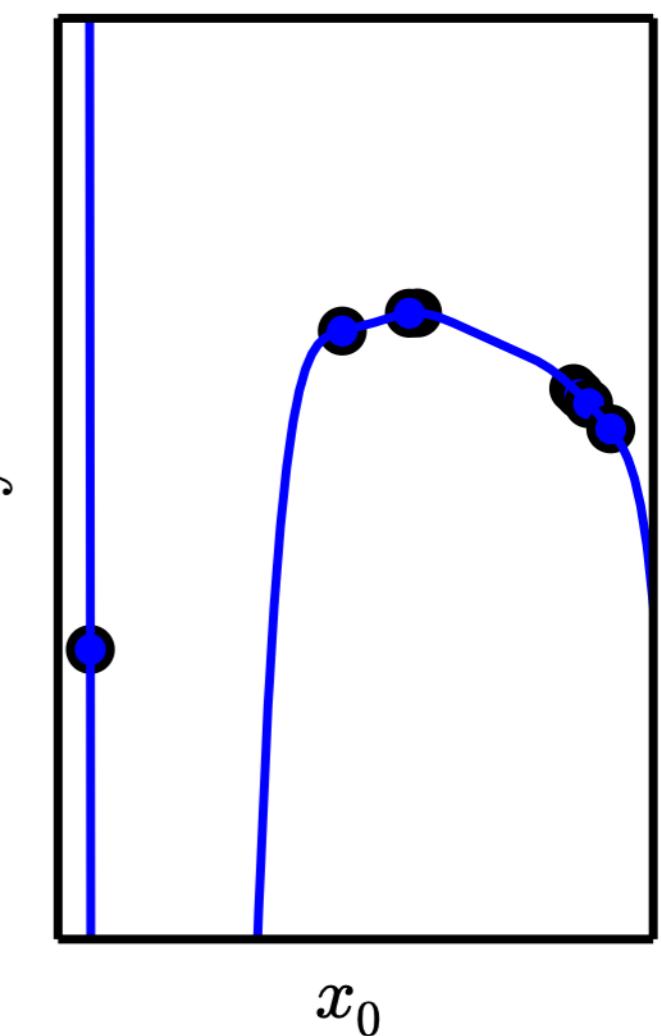


Figure 5.5

L2-norm regularization with Linear Regression

↳ Ridge

(Review)
 $\ell-1$ norm \rightarrow Lasso

$$\hat{y}^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)}$$

$(\beta = \mathbf{w})$

Design Matrix / Data Matrix $\mathbf{X} \in \mathbb{R}^{n \times m}$

n - # samples
 m - # features

$$\min_{\mathbf{w}} \frac{1}{2} \|\hat{y} - y^*\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

$$\frac{1}{2} \|\mathbf{X}\mathbf{w} - \hat{y}^*\|_2^2 + \lambda \mathbf{w}^T \mathbf{w}$$

$$\frac{1}{2} (\mathbf{X}\mathbf{w} - \hat{y}^*)^T (\mathbf{X}\mathbf{w} - \hat{y}^*) + \lambda \mathbf{w}^T \mathbf{w}$$

$$\frac{1}{2} (\underline{\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}} - 2 \underline{\mathbf{w}^T \mathbf{X}^T \hat{y}^*} + \cancel{\hat{y}^T \hat{y}^*}) + \underline{\lambda \mathbf{w}^T \mathbf{w}}$$

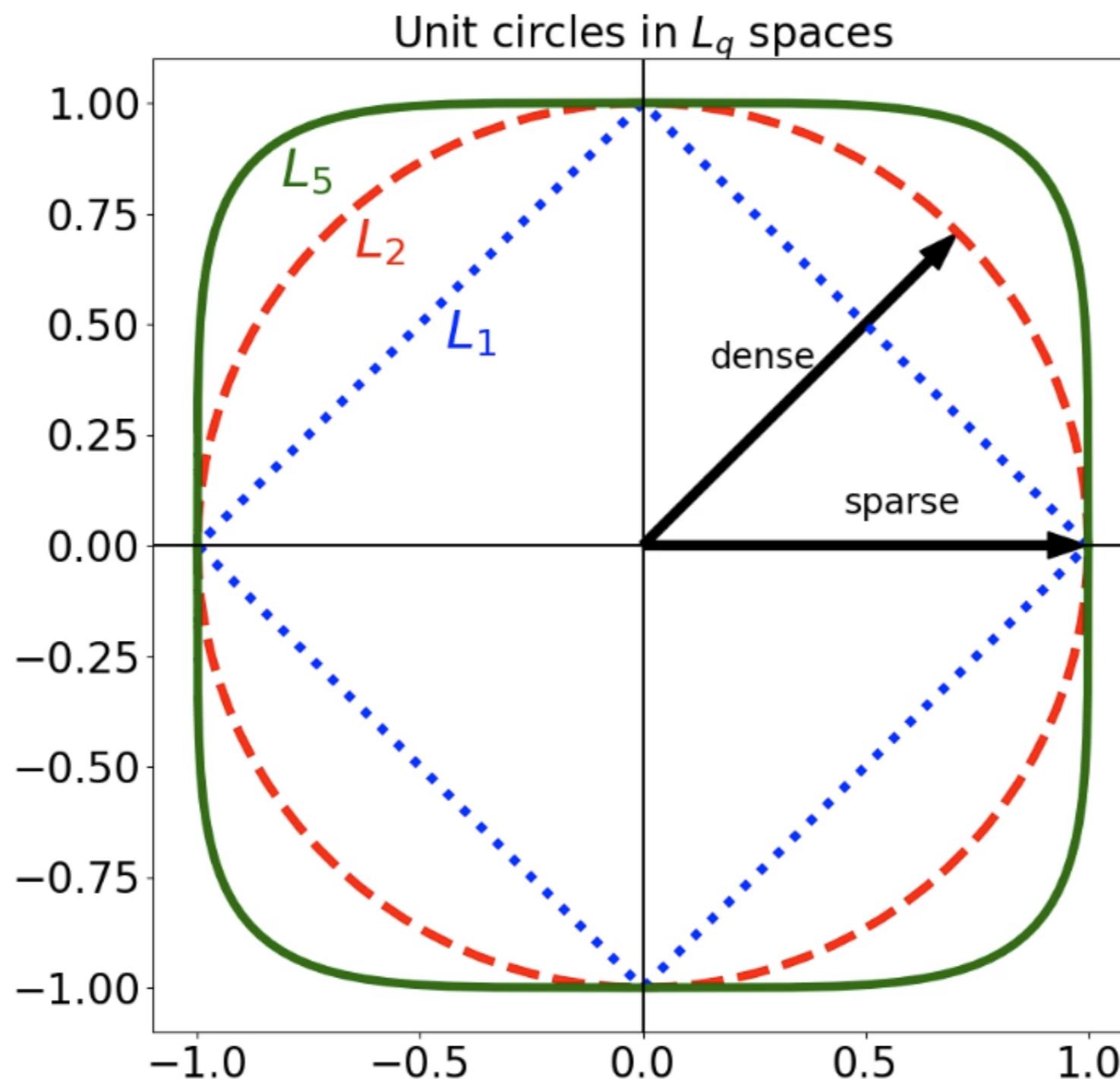
$$\nabla_{\mathbf{w}} = \underline{\mathbf{X}^T \mathbf{X} \mathbf{w}} - \underline{\mathbf{X}^T \hat{y}^*} + \underline{2\lambda \mathbf{w}}$$

$$= (\mathbf{X}^T \mathbf{X} + 2\lambda \mathbf{I}) \mathbf{w} - \mathbf{X}^T \hat{y}^* = \mathbf{0}$$

Solve for \mathbf{w} $\rightarrow \mathbf{w} = (\mathbf{X}^T \mathbf{X} + 2\lambda \mathbf{I})^{-1} \mathbf{X}^T \hat{y}^*$

Parameter Norm Penalties

Norm



Training

- Define loss (objective) function ✓
- Compute gradient of loss w.r.t. parameters
- Update parameters by moving “downhill” in the direction of the negative gradient ✓
- If loss stops decreasing (or other convergence criteria satisfied), you are done! ✓

$f(x + \Delta x)$

~~$f(x)$~~

$\nabla_x f(x)$

Computing Gradients

- ✓ 1. **Analytically** (with equations for partial derivatives)
- ✓ 2. **Numerically** (with repeated, perturbed forward solves) $\nabla_x \sim f(x + \Delta x) - f(x)$
- * 3. **With Automatic Differentiation** (built-in to tensorflow, pytorch, etc.)

Gradient Descent

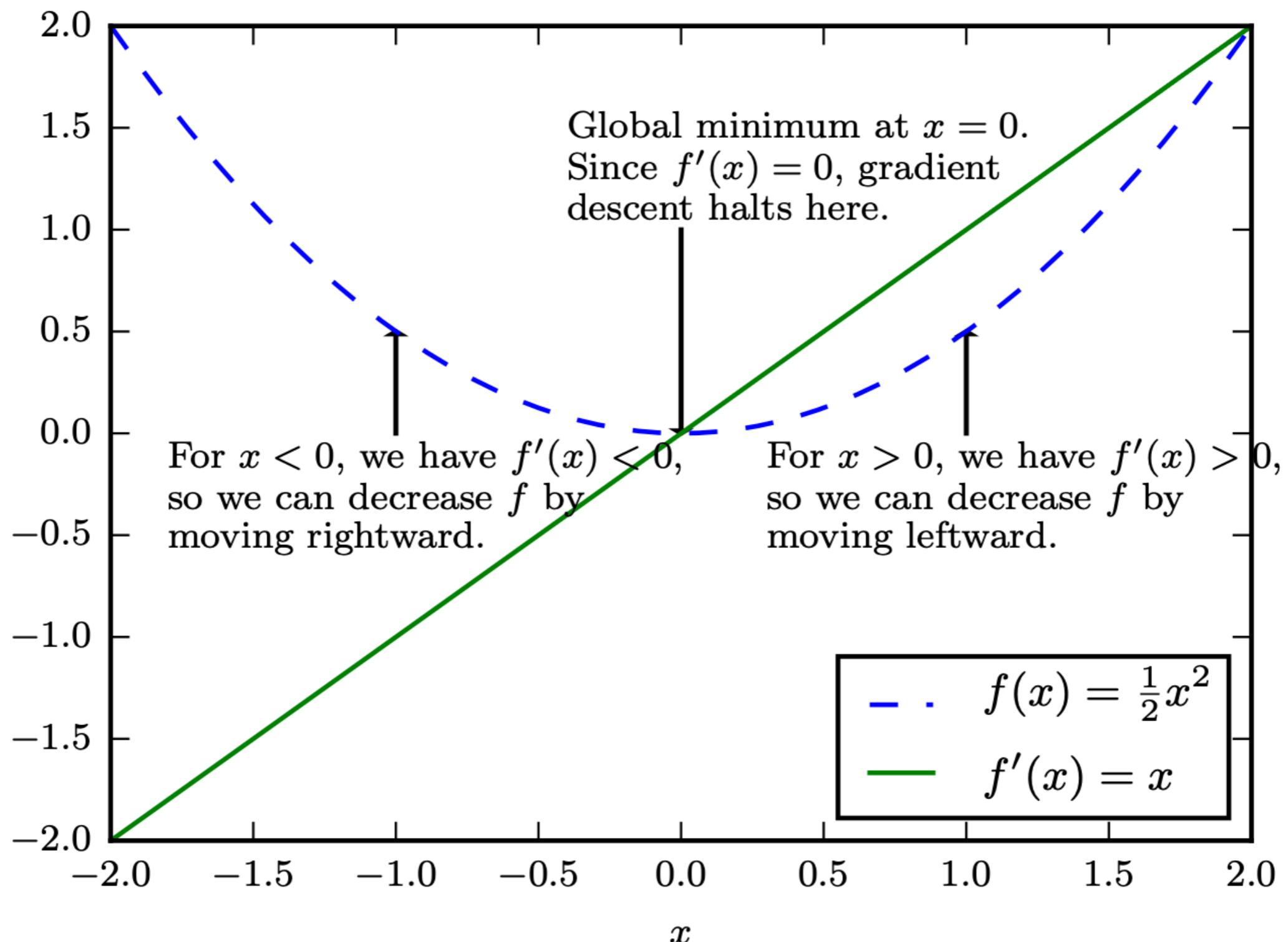


Figure 4.1

Gradient Descent with Poor Conditioning

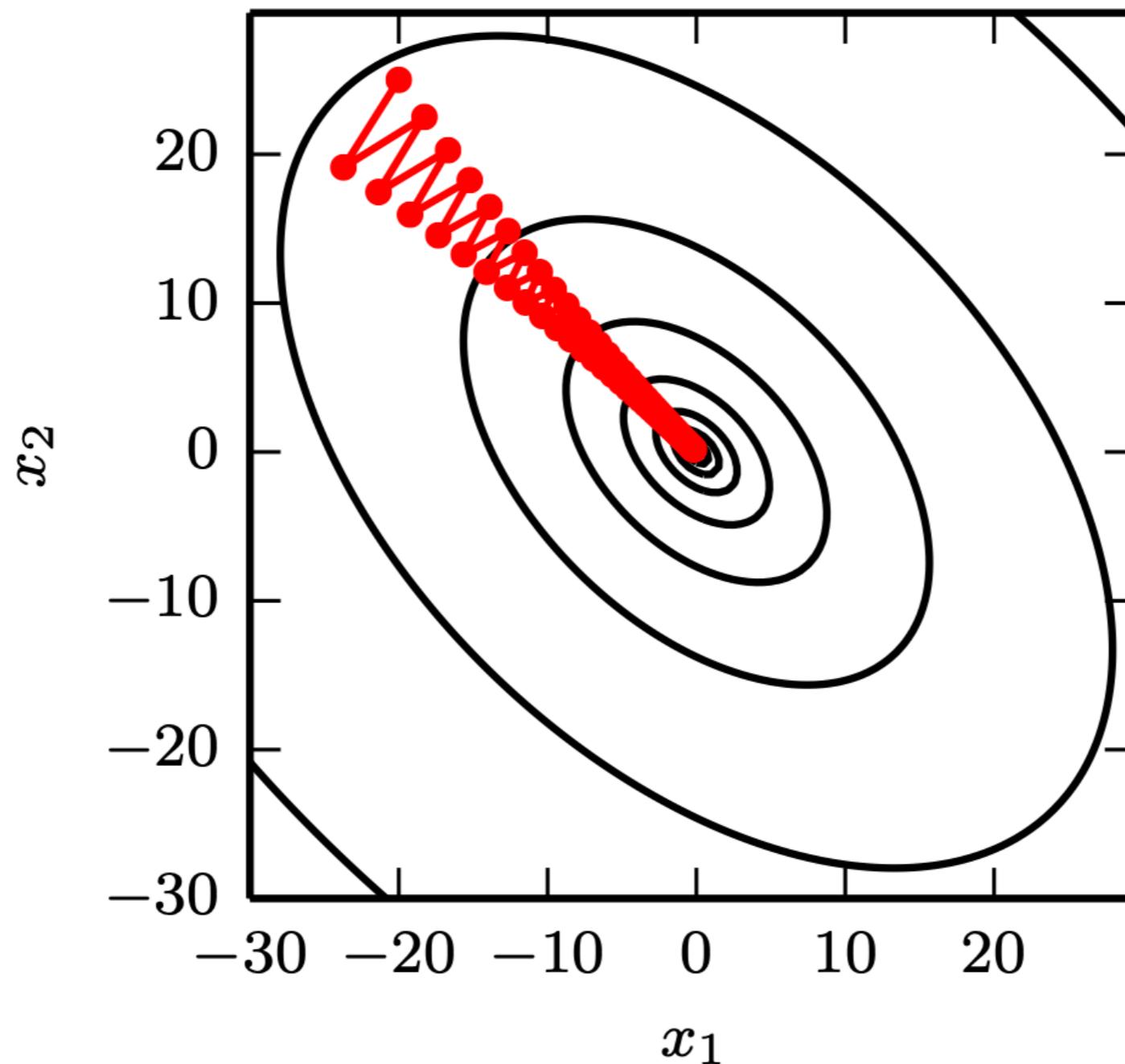
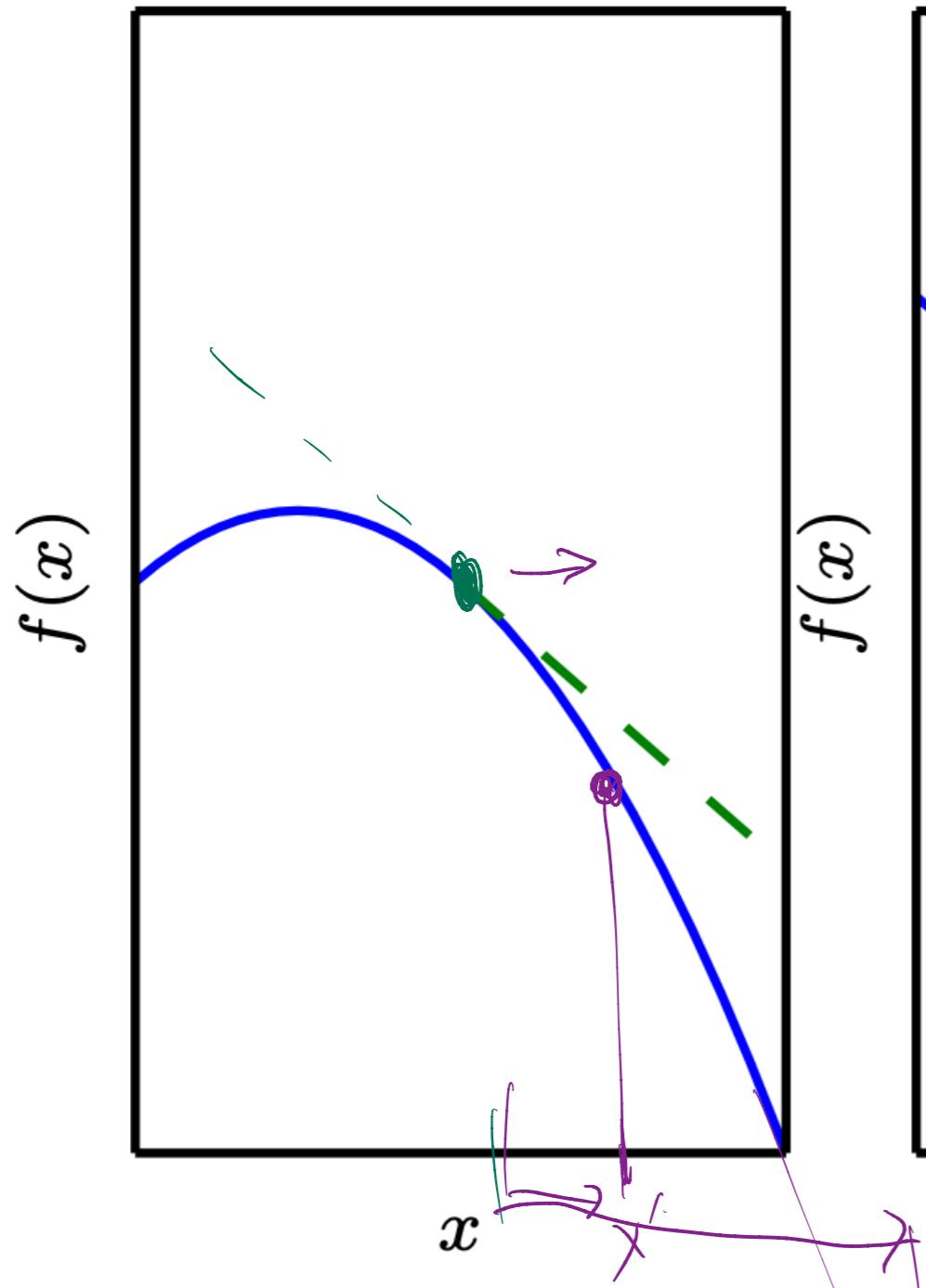


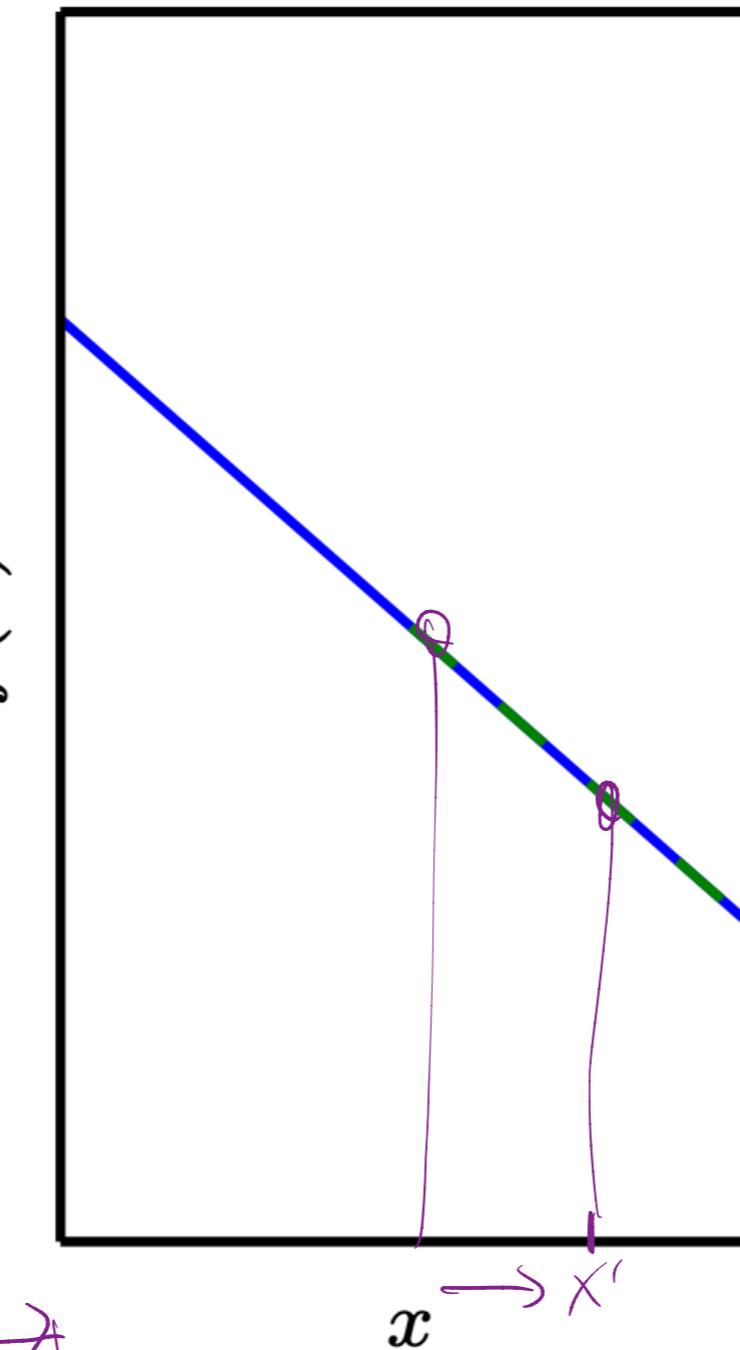
Figure 4.6

Curvature

Negative curvature



No curvature



Positive curvature

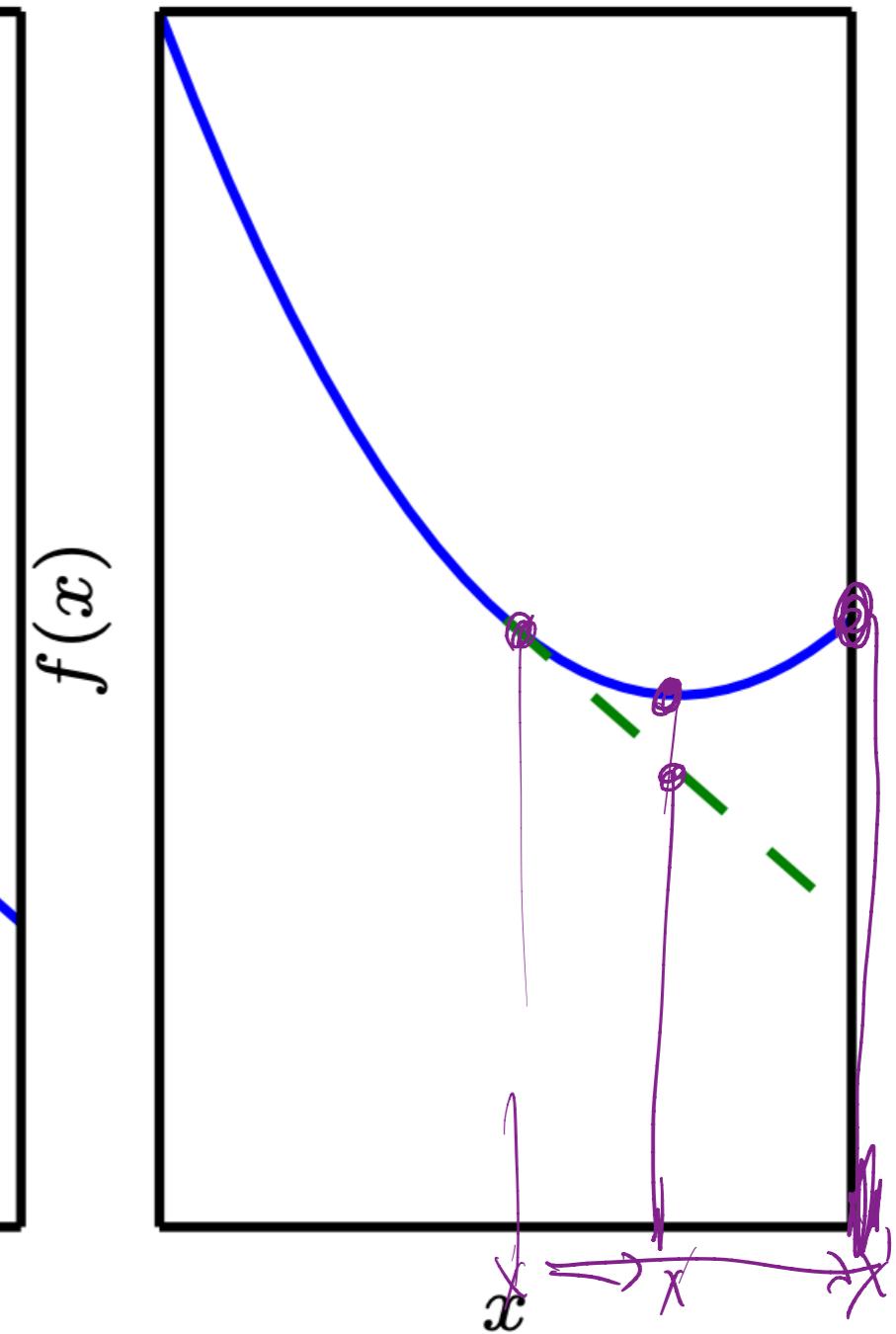


Figure 4.4

Newton's Method $\xrightarrow{\text{+}} \text{Raphson}$ (trade off b/w 1st / 2nd order)

Taylor series of loss

$$f(x)$$

$$f(x + \Delta x) = f(x) + f' \Delta x + f'' \Delta x^2 + \dots$$

1st order

2nd order

1st derivative info
Gradient

2nd derivative information
"Hessian"

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)}) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H}(f)(\mathbf{x}^{(0)})(\mathbf{x} - \mathbf{x}^{(0)}).$$

$$\mathbf{x}^* = \mathbf{x}^{(0)} - \mathbf{H}(f)(\mathbf{x}^{(0)})^{-1} \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)}). \quad (4.11)$$

Predicting optimal step size using Taylor series

$$f(\mathbf{x}^{(0)} - \epsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \epsilon \mathbf{g}^\top \mathbf{g} + \frac{1}{2} \epsilon^2 \mathbf{g}^\top H \mathbf{g}. \quad (4.9)$$

$$\epsilon^* = \frac{\mathbf{g}^\top \mathbf{g}}{\mathbf{g}^\top H \mathbf{g}}. \quad (4.10)$$

Big gradients speed you up

Big eigenvalues slow you
down if you align with their
eigenvectors

Challenges in Neural Network Optimization

III-Conditioning Local Minima

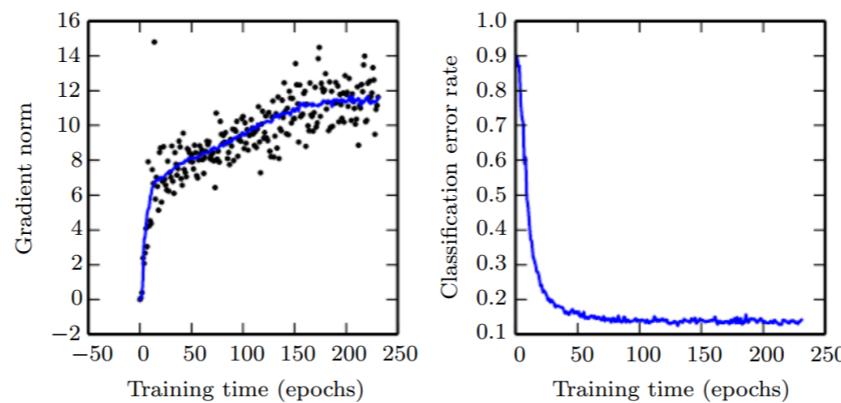


Fig 8.1 Goodfellow et al.
2016

Plateaus, Saddle Points and Flat regions

Cliffs and Exploding gradients

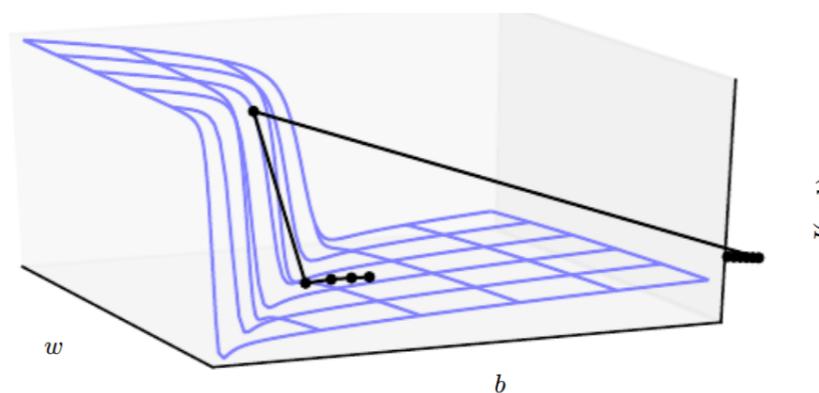


Fig 8.3 Goodfellow et
al. 2016

Inexact Gradients

Long Term Dependencies

Gradient Descent with Mini-batches

Larger batch size → more accurate estimate of gradient

Memory required scales with batch size

Batch sizes of power of 2 (32, 64, 128, 256)

For stability use a lower learning rate if batch size is smaller

Select minibatches randomly so that the samples are independent

Stochastic Gradient Descent

Algorithm 7: Stochastic gradient descent (SGD) at training iteration k

Require: Learning rate ϵ_k

Require: Initial parameter θ

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set

$\{x(1), \dots, x(m)\}$ with corresponding targets $y(i)$;

 Compute gradient estimate: $\hat{g} \leftarrow \frac{1}{m} \nabla \sum_i L(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$;

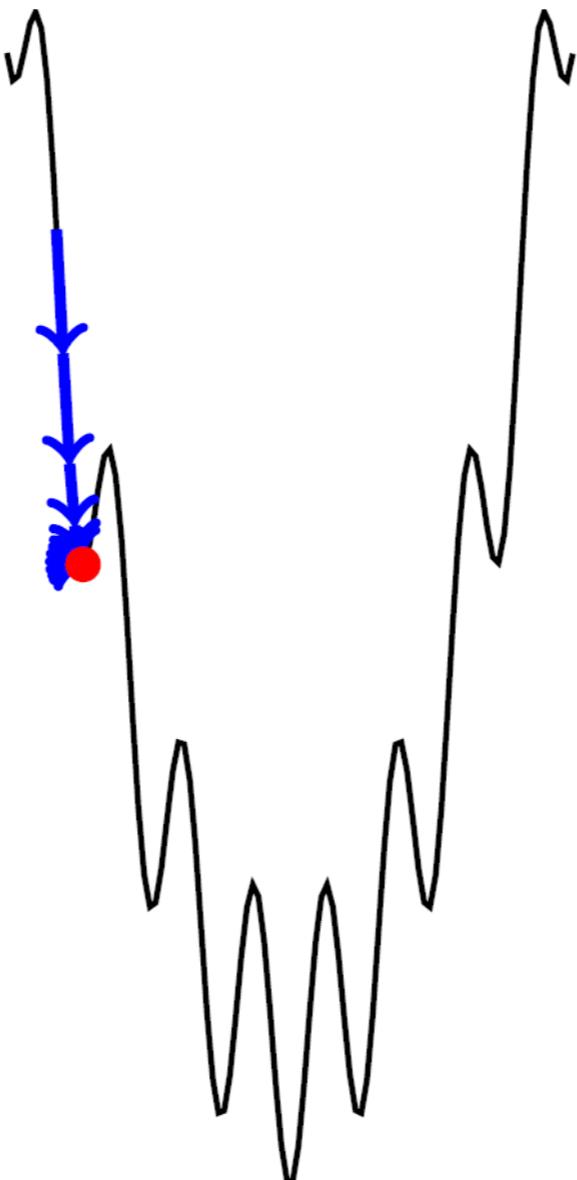
 Apply update: $\theta \leftarrow \theta - \epsilon \hat{g}$;

end

m is subset of training samples

Stochastic Gradient Descent

Learning rate: too low

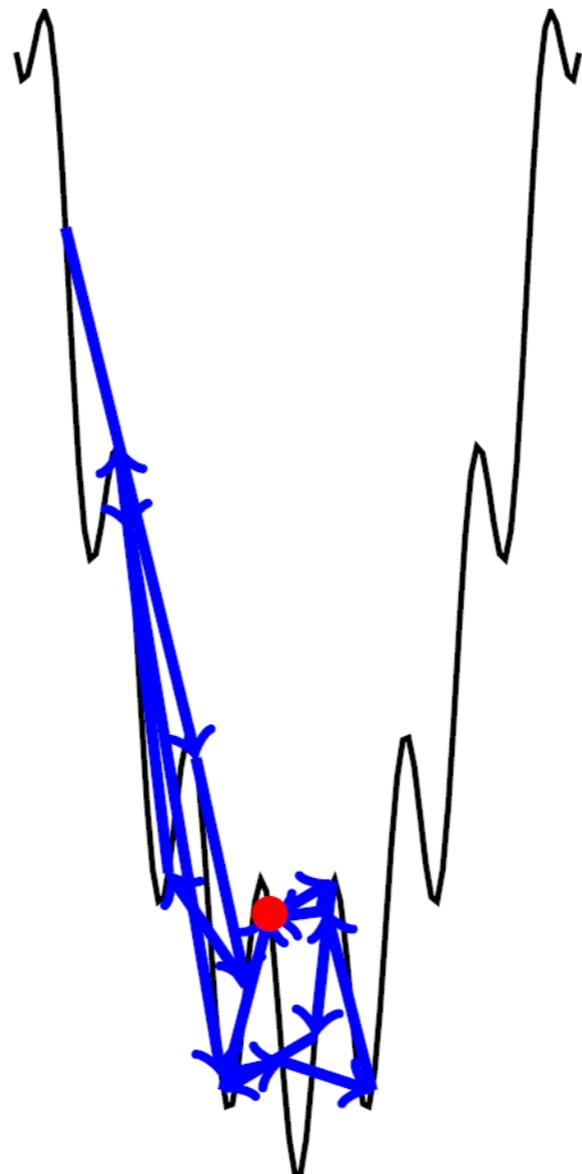


$$\nabla J = -0.0538 \quad \mathbf{w} = -2.92369$$

$$n = 14 \quad \eta = 0.005$$

Stochastic Gradient Descent

Learning rate: too high

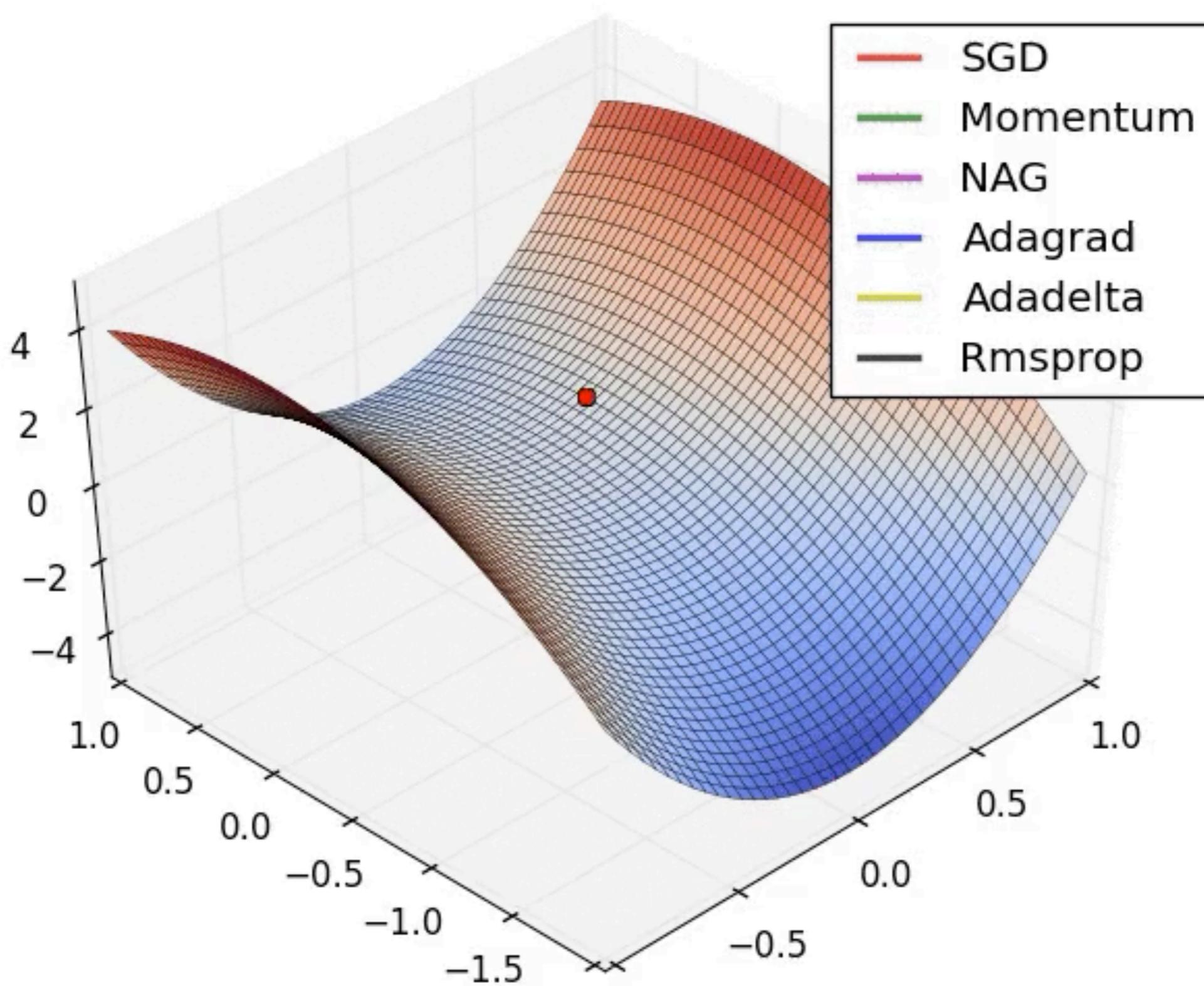


$$\nabla J = -8.28964 \quad \mathbf{w} = -0.39825$$

$$n = 14 \quad \eta = 0.1$$

Rate scheduling
built-in, e.g. Adam

Visualize Adaptive Learning Rates



<https://i.imgur.com/2dKCQHh.gif>

Dataset Augmentation

A classifier should be invariant to a wide variety of transformations

Generate new (x, y) pairs by transforming the x inputs in our training set:

- translation

- rotation

- scaling

- cropping

Not always possible, e.g., out-of-plane rotation

General idea: noise injection (at the inputs x)

Injecting Noise at the Output Targets

Label Smoothing:

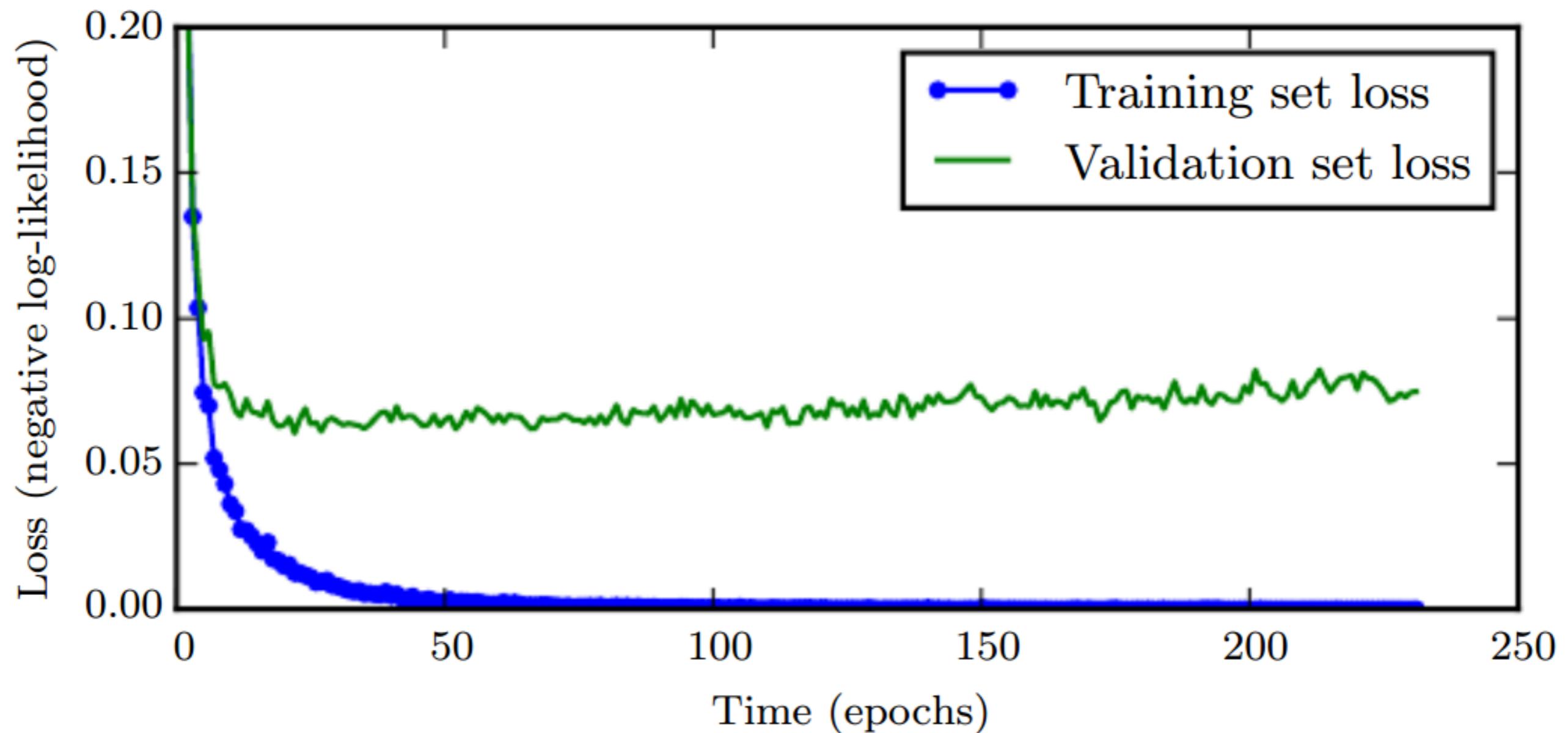
Generally the target is

$$y = [0, 0, 0, \textcolor{blue}{1}, 0]$$

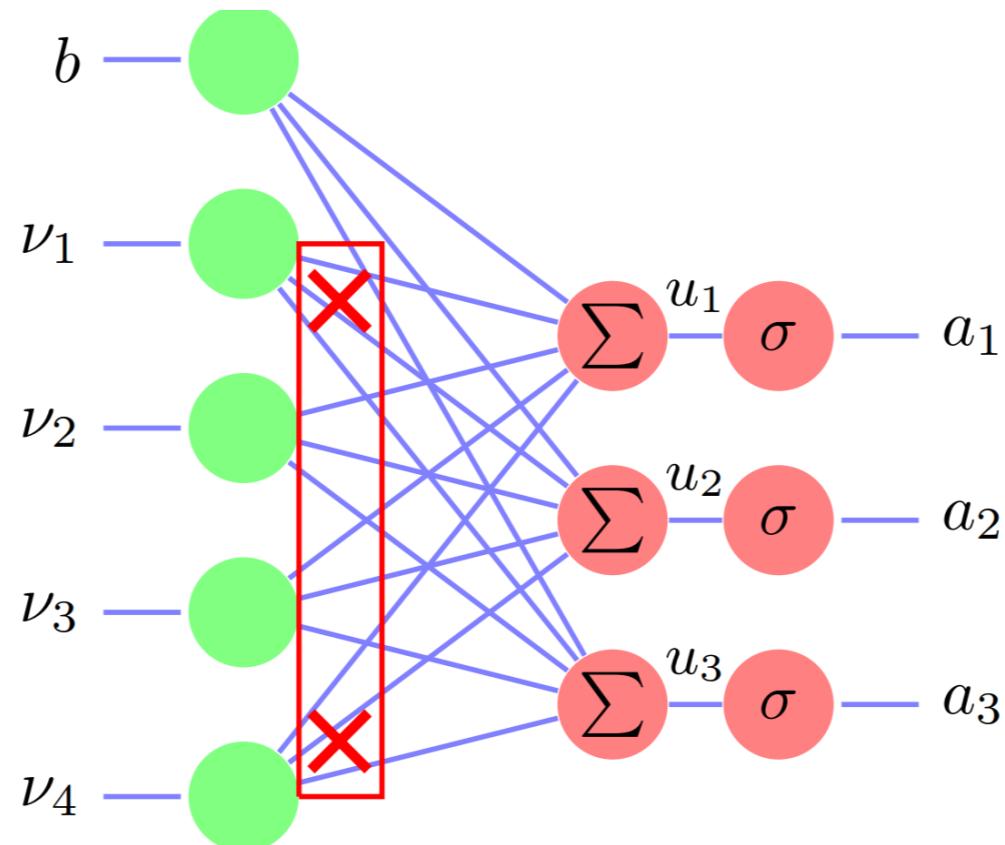
Use instead:

$$y = [\frac{\epsilon}{k-1}, \frac{\epsilon}{k-1}, \frac{\epsilon}{k-1}, \textcolor{blue}{1-\epsilon}, \frac{\epsilon}{k-1}]$$

Early Stopping



Dropout



Dropout

