

Artificial Intelligence and Machine Learning

Applications

Lecture Outline

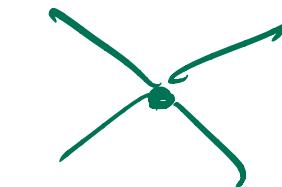
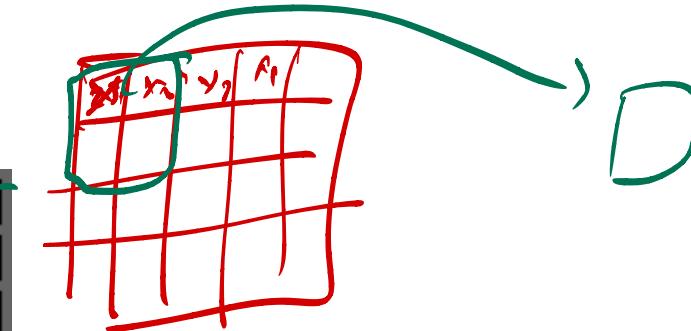
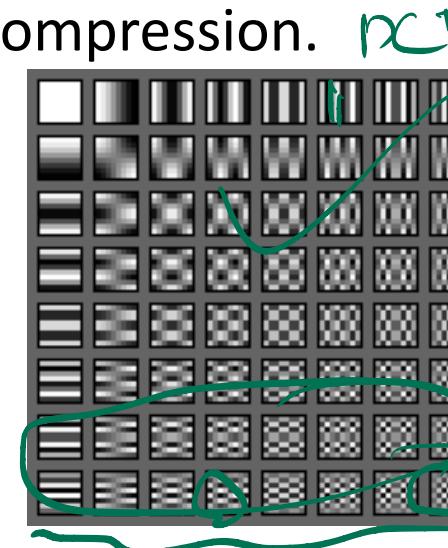
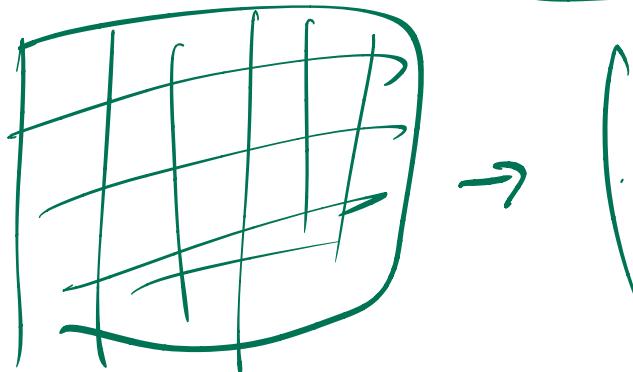
- Linear Regression Tutorial & Playground
- Image Compression
- Auto-Encoders
- Optimizers
- Logistic Regression Interactive Demo
- Neural Network Regression Tutorial & Playground

Image Compression with Linear Regression

10

$$X \rightarrow f(X) = X' \Rightarrow \text{CompFunc}(X') \rightarrow X_{\text{low}}$$

- An image can be transformed into the frequency domain and represented as a combination of some basic components.
- Cosine Basis Functions and Discrete Cosine Transforms (DCT) can enable this.
- The human eye is most sensitive to low frequencies. Therefore, most of the high frequencies can be ignored.
- Principal behind JPEG Image Compression.



$$f^k(X_{\text{low}})$$

$$X_{\text{comp}}$$

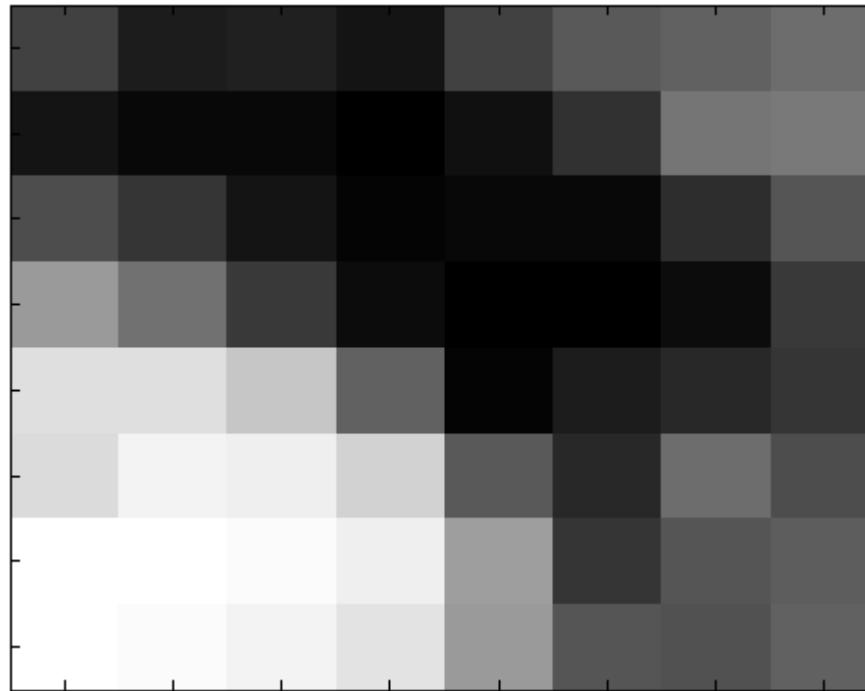
Image Compression with Linear Regression

- Using linear regression, we can learn the weight of each of these cosine basis.
- We only need to store the weights of the basis frequencies.
- The image is reconstructed using these weights.

ـ ـ ـ
ـ ـ ـ
ـ ـ ـ

Image Compression

8 x 8 Pixels



Image



Image Compression

- Gray-Scale Example:
- Value Range 0 (black) --- 255 (white)

63	33	36	28	63	81	86	98
27	18	17	11	22	48	104	108
72	52	28	15	17	16	47	77
132	100	56	19	10	9	21	55
187	186	166	88	13	34	43	51
184	203	199	177	82	44	97	73
211	214	208	198	134	52	78	83
211	210	203	191	133	79	74	86

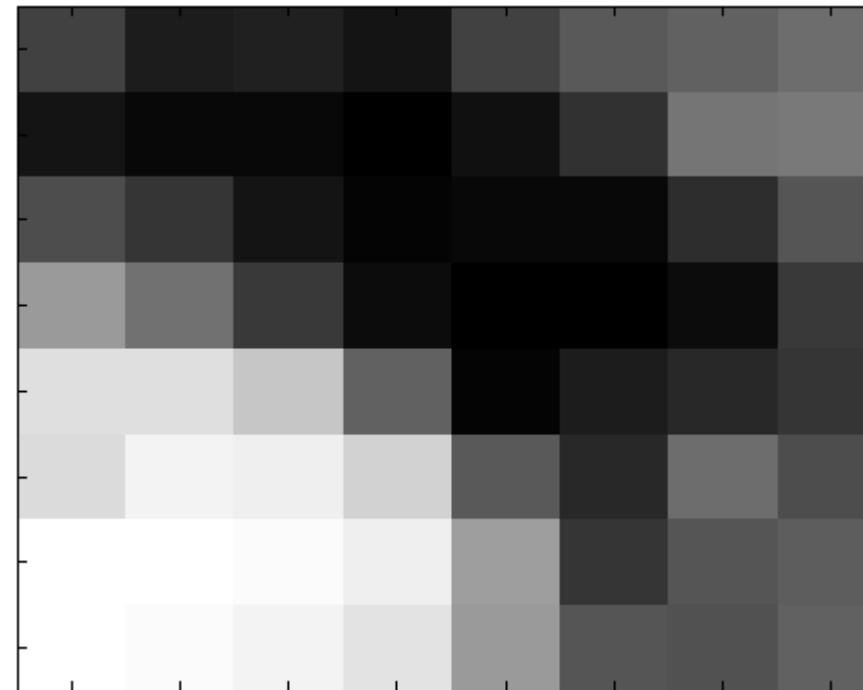


Image Compression

- **2D-DCT of matrix**

Numbers are coefficients
of polynomial

```
-304 210 104 -69 10 20 -12 7
-327 -260 67 70 -10 -15 21 8
 93 -84 -66 16 24 -2 -5 9
 89  33 -19 -20 -26 21 -3 0
 -9  42 18 27 -7 -17 29 -7
 -5  15 -10 17 32 -15 -4 7
 10   3 -12 -1  2  3 -2 -3
 12  30  0 -3 -3 -6 12 -1
```

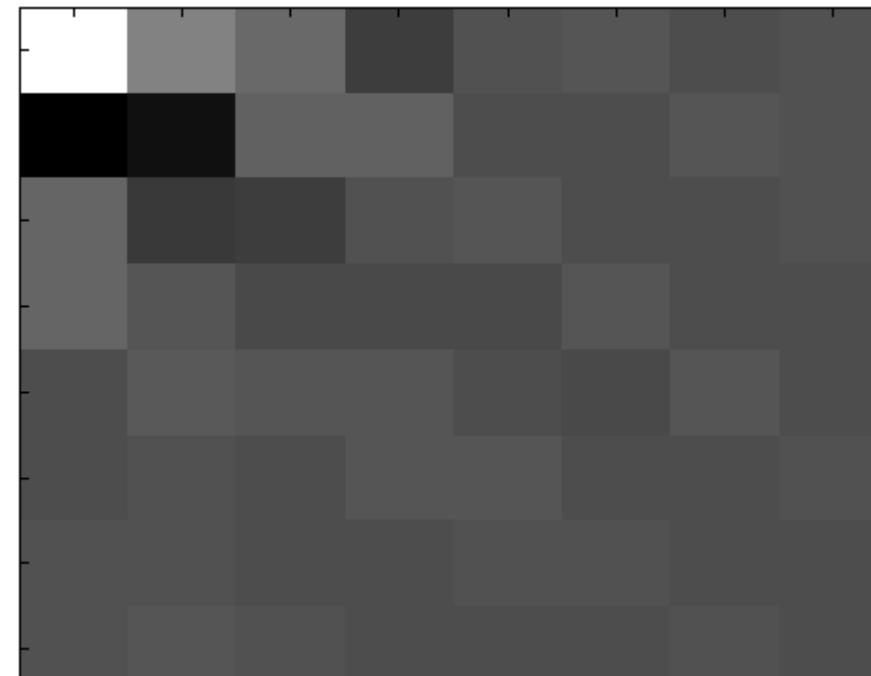


Image Compression

- Cut the least significant components

-304 210 104 -69 10 20 -12 0
-327 -260 67 70 -10 -15 0 0
93 -84 -66 16 24 0 0 0
89 33 -19 -20 0 0 0 0
-9 42 18 0 0 0 0 0
-5 15 0 0 0 0 0 0
10 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0

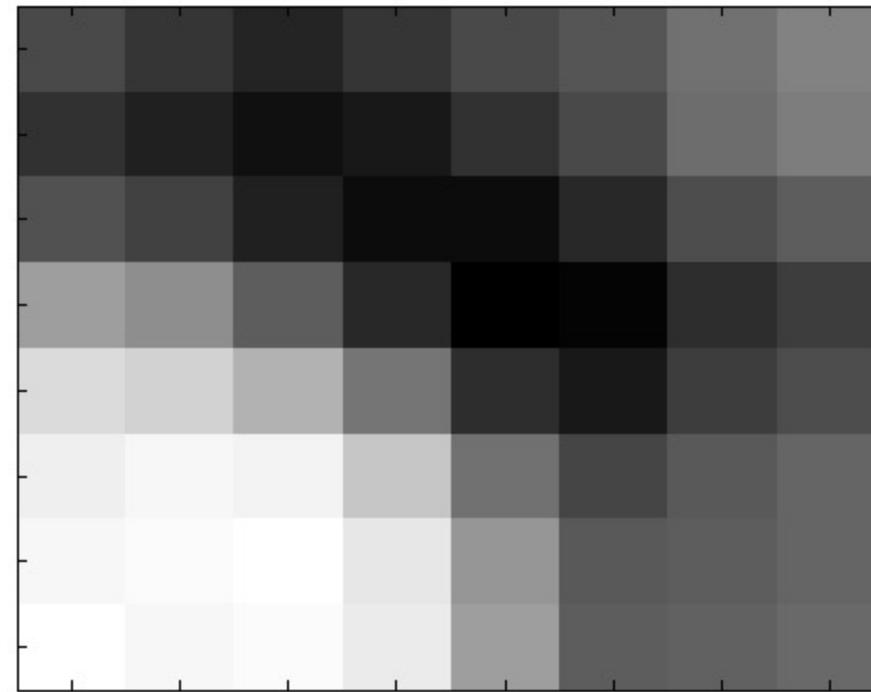


As you can see, we save a little over half the original memory.

Reconstructing the Image

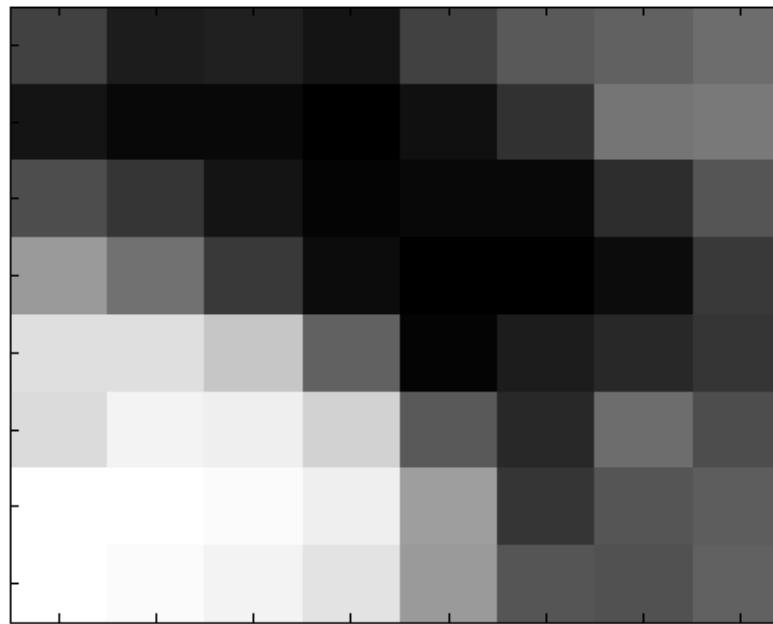
- New Matrix and Compressed Image

55	41	27	39	56	69	92	106
35	22	7	16	35	59	88	101
65	49	21	5	6	28	62	73
130	114	75	28	-7	-1	33	46
180	175	148	95	33	16	45	59
200	206	203	165	92	55	71	82
205	207	214	193	121	70	75	83
214	205	209	196	129	75	78	85



Can You Tell the Difference?

Original



Compressed

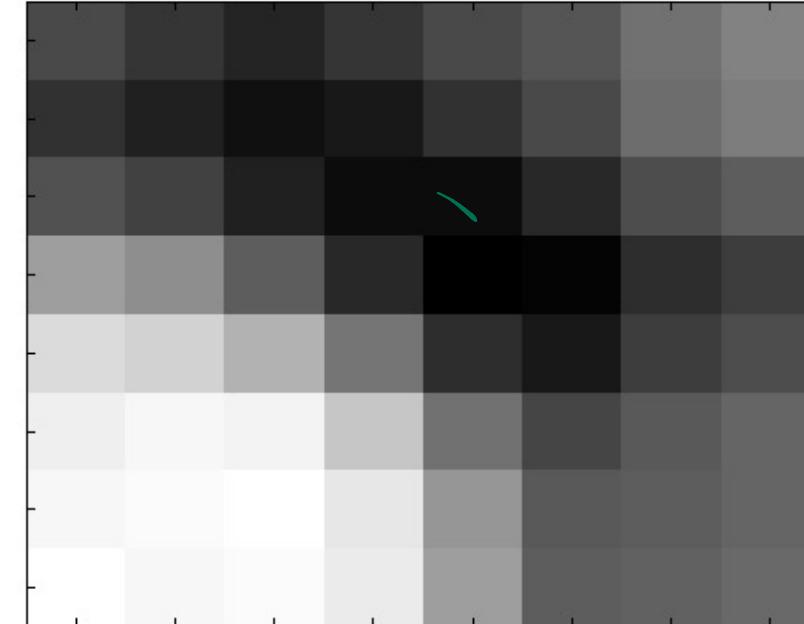


Image Compression

Original



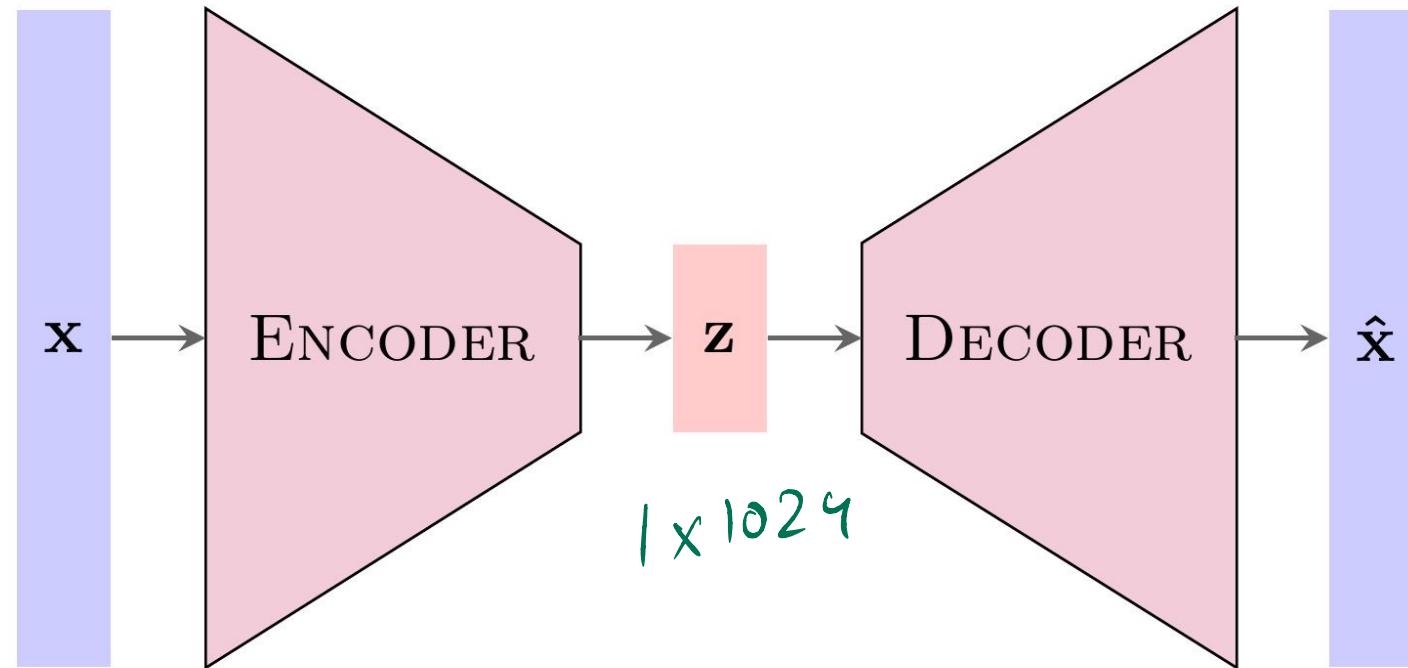
Compressed



Applications – AutoEncoders

- Autoencoders are a type of neural networks where the input is also the output.
- They come under unsupervised learning and there are no labels involved.
- An autoencoder consists of two parts: encoder and decoder.
- The idea here is that you take a higher dimensional input, project it into a lower dimensional space and then project it back into the input space.

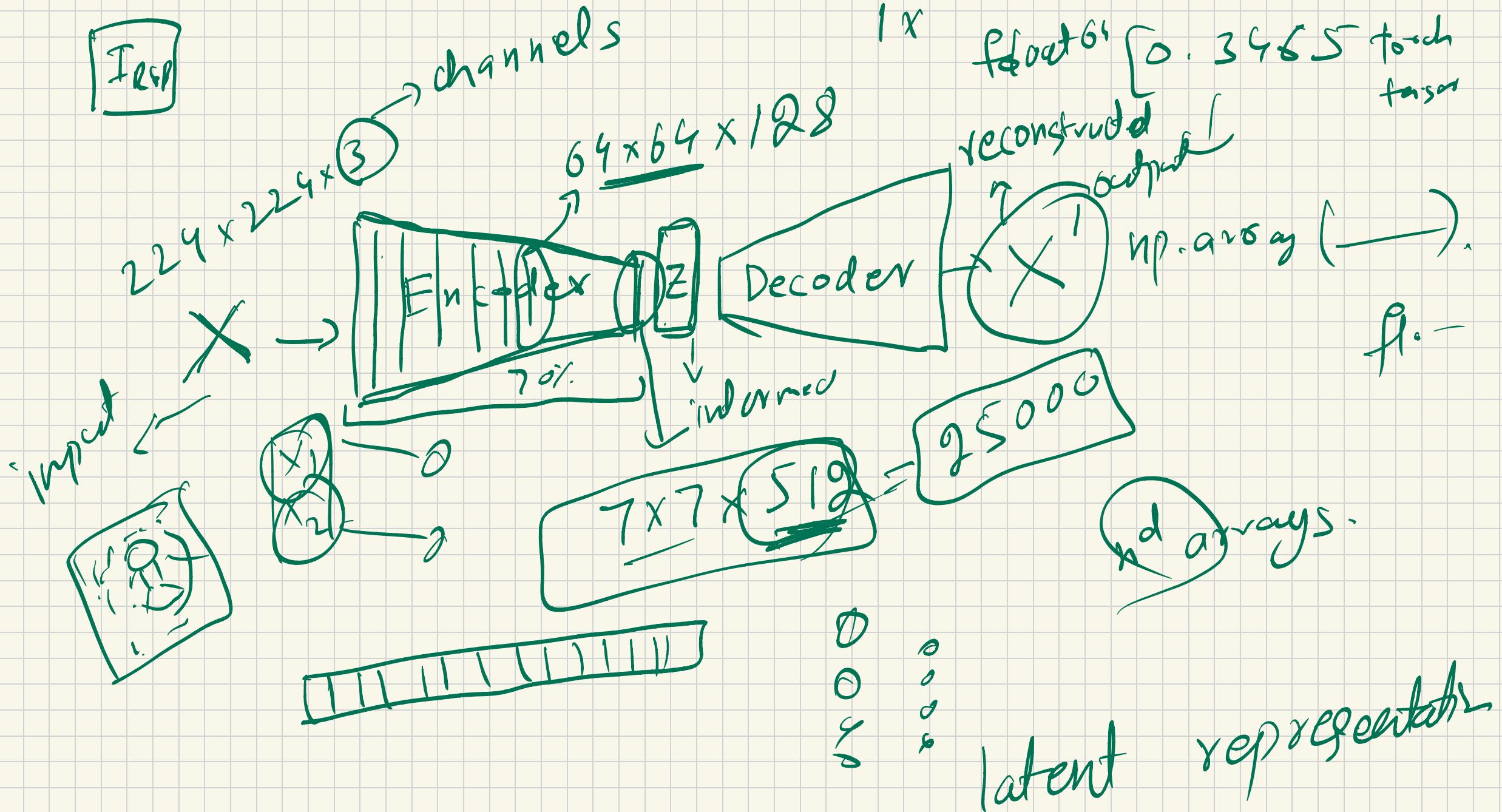
Applications – AutoEncoders

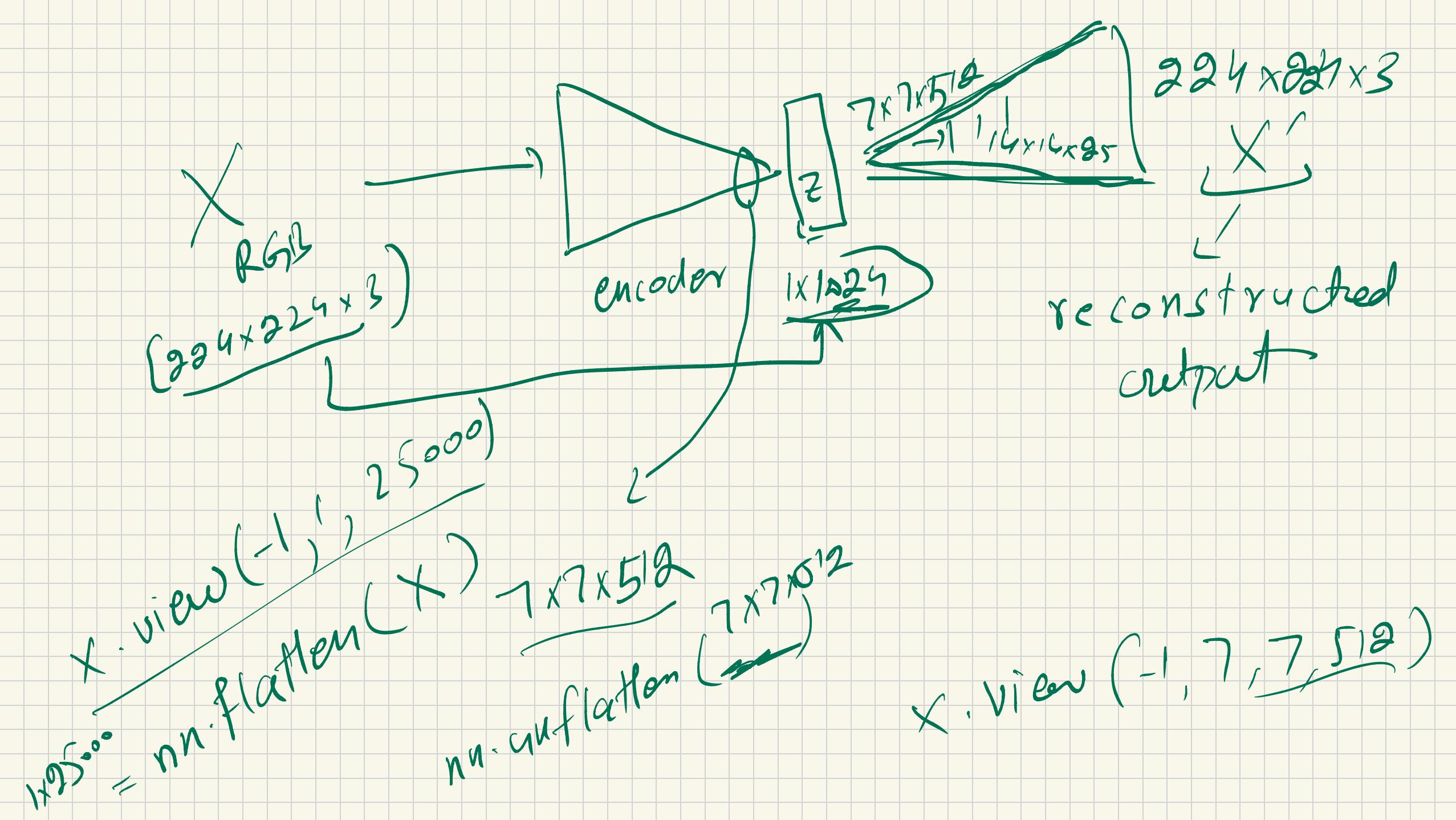


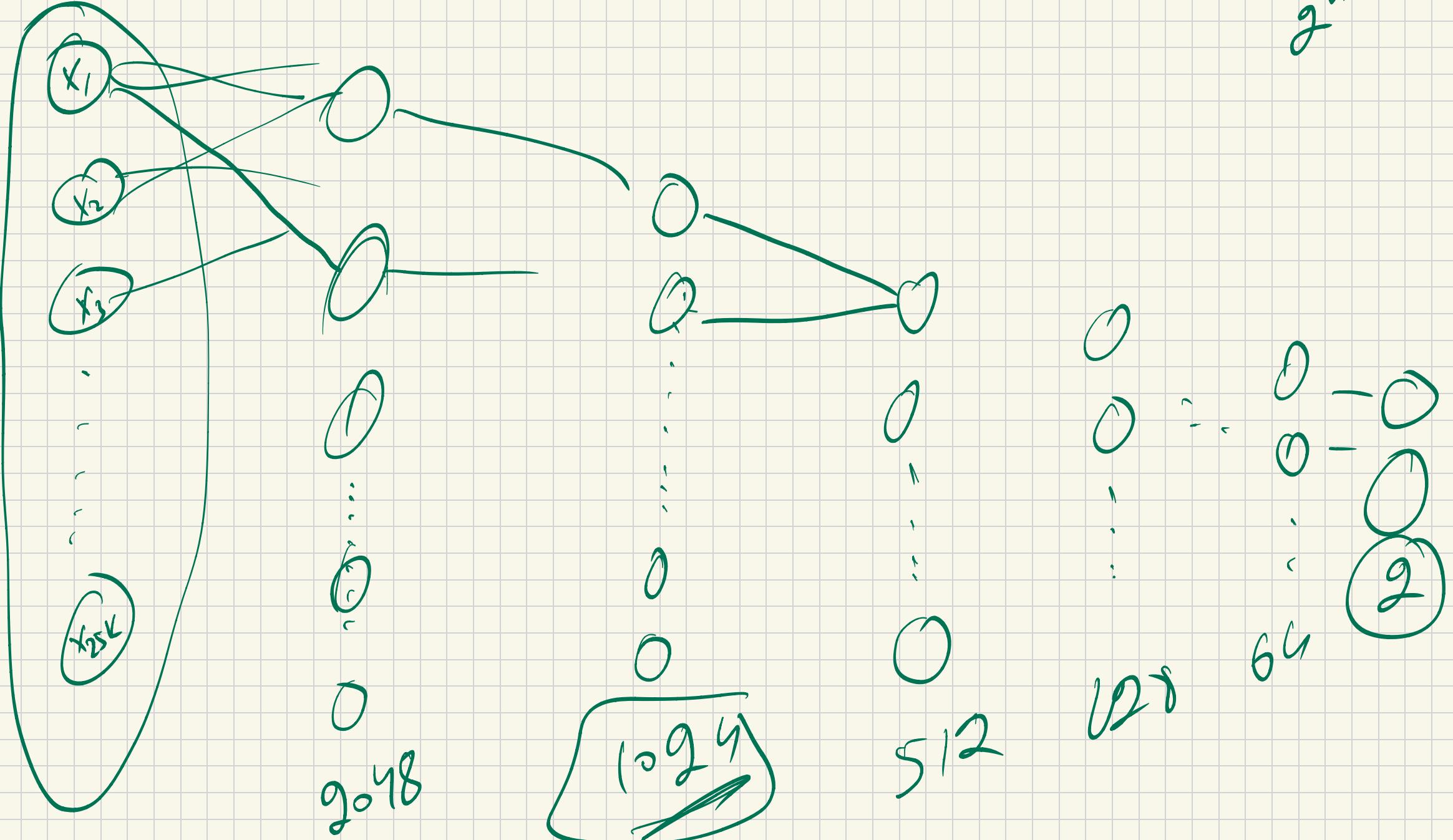
Applications – AutoEncoders

- The autoencoder model tries to minimize the reconstruction error (RE).
- Typically, mean squared is used as the loss function for autoencoders.
- The objective is to minimize the following:

$$L(x, \hat{x}) = \frac{1}{N} \sum_{i=1}^N ||x_i - \hat{x}_i||^2$$

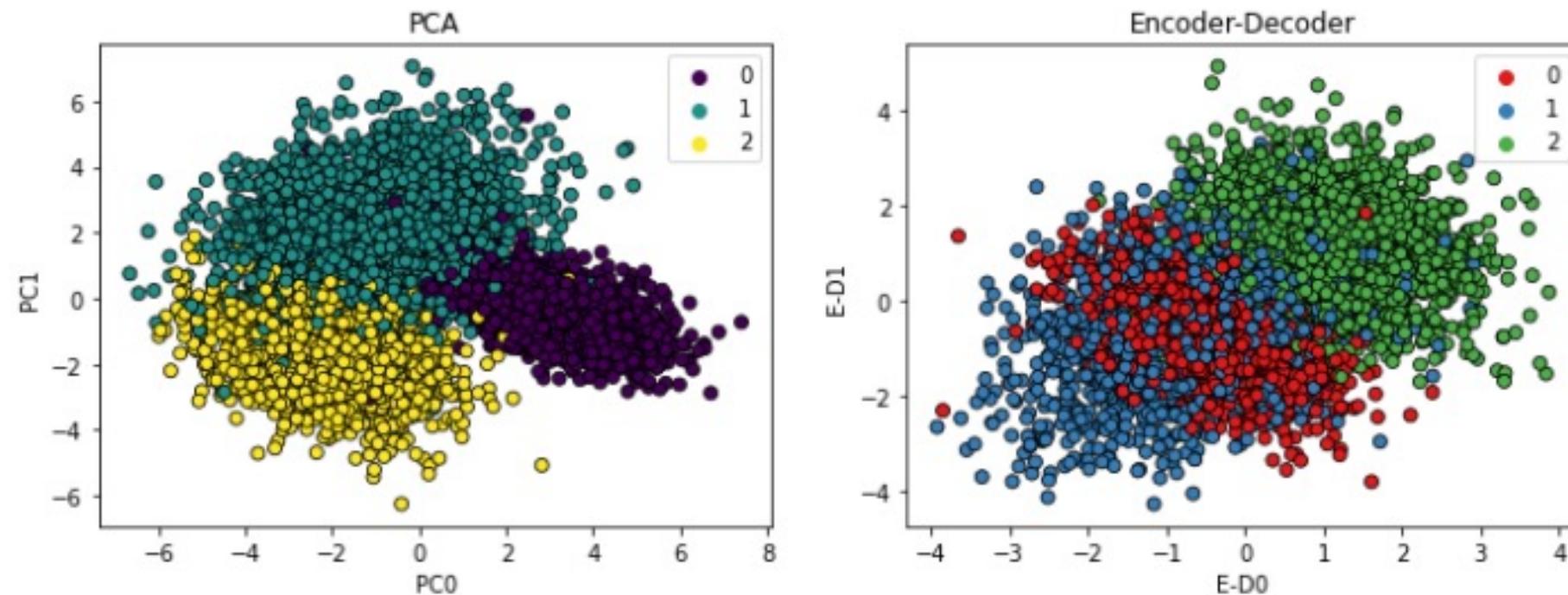






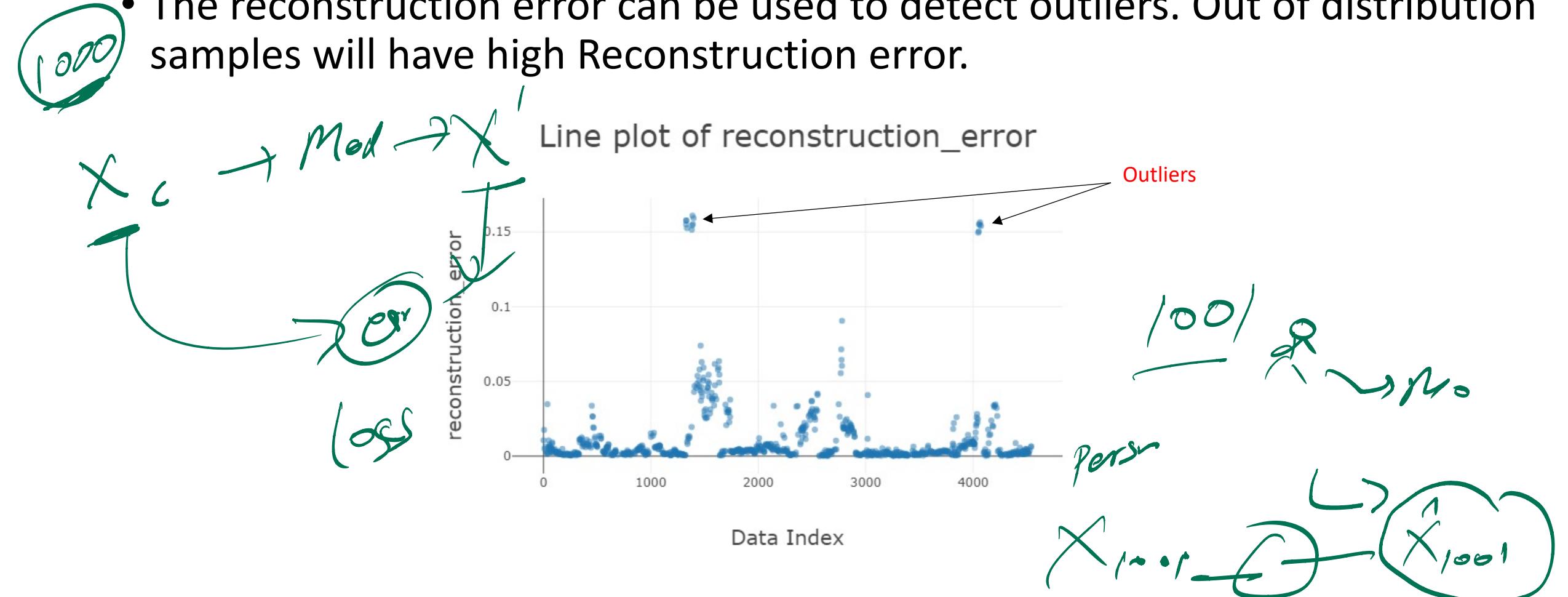
AutoEncoders for Dimensionality Reduction

- The encoder output can be used for dimensionality reduction

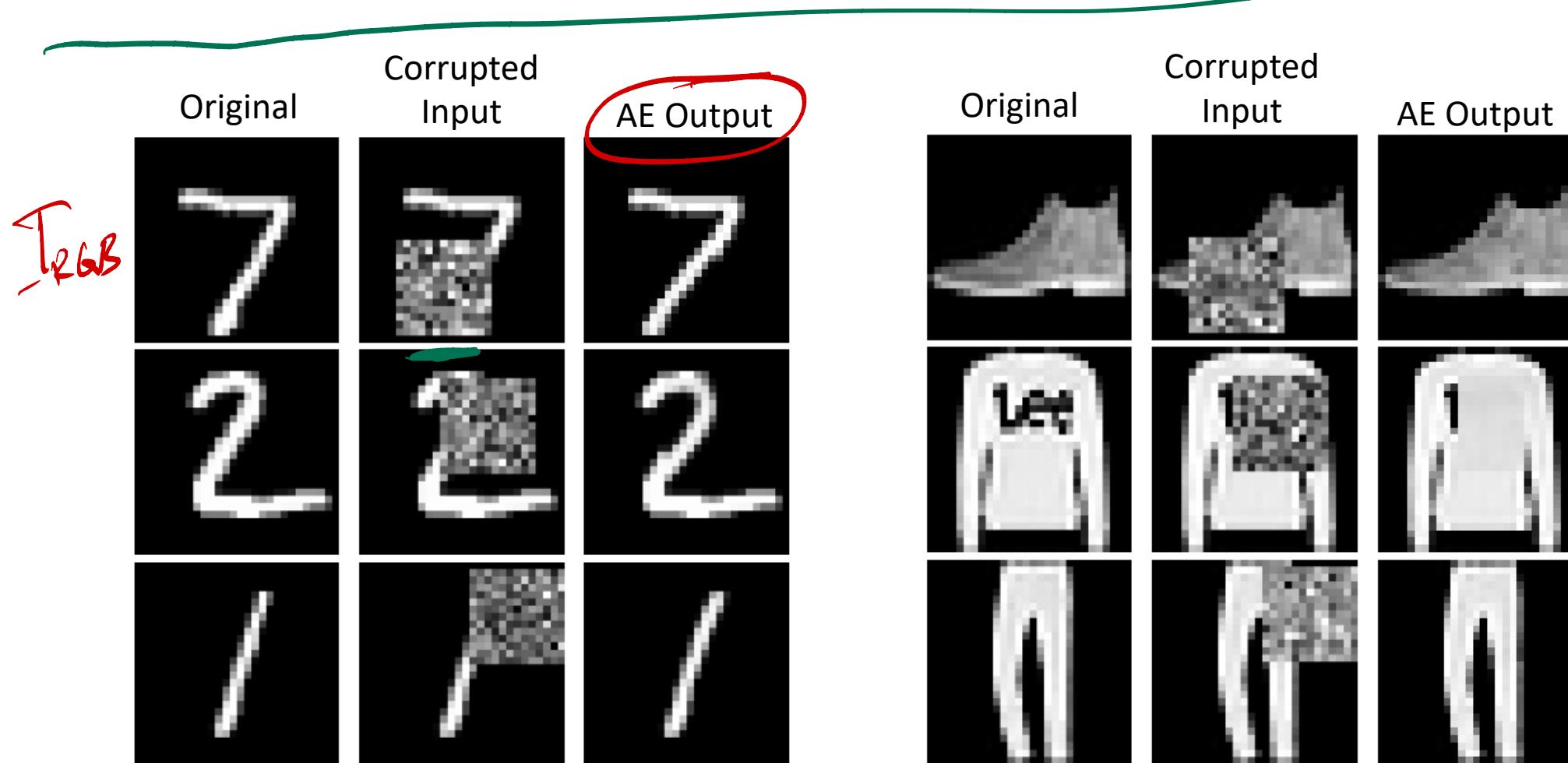


AutoEncoders for Outlier Detection

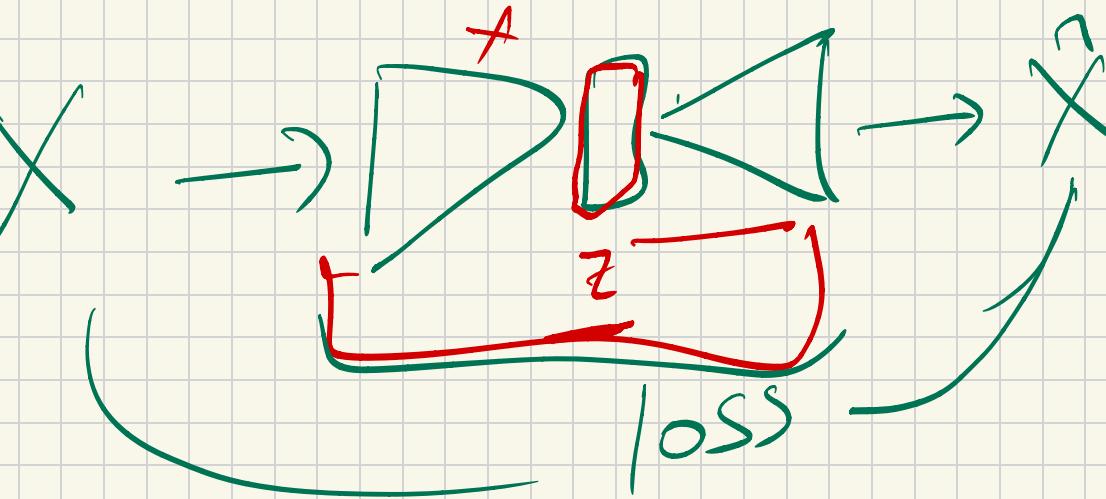
- The reconstruction error can be used to detect outliers. Out of distribution samples will have high Reconstruction error.



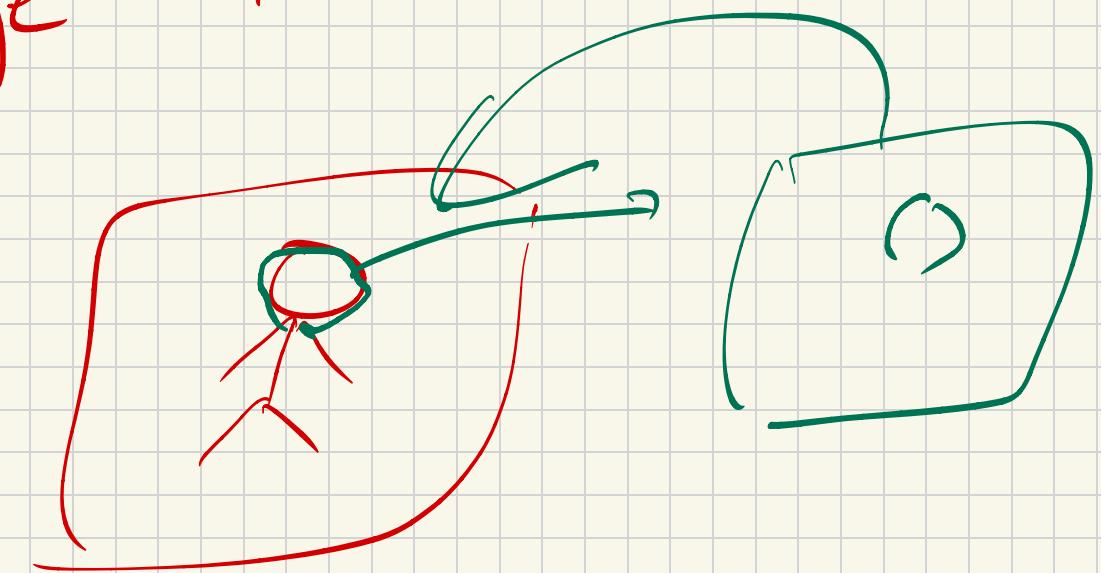
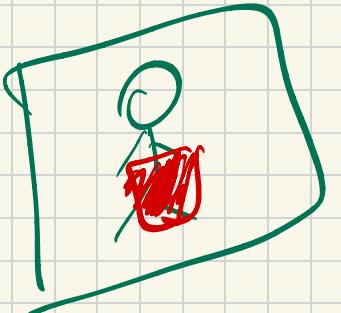
AutoEncoders for Image Completion



↑
RGB



Generative
Image



ML Algorithms Perspectives:

- We can look into ML algorithms from two perspective:
 - **Loss Minimization** Problem (like what we did).
 - **Probability Maximization** Problem (using Maximum Likelihood Estimation).
- Both should result in the same solution.

Probalistic Interpretation of Linear Regression and MLE

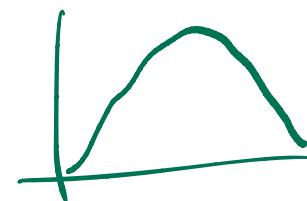
- We can also look at the probalistic Interpretation of Linear Regression.
- Keeping everything else same as the previous formulation

$$y_i = \mathbf{x}_i^T \boldsymbol{\theta} + \epsilon_i$$

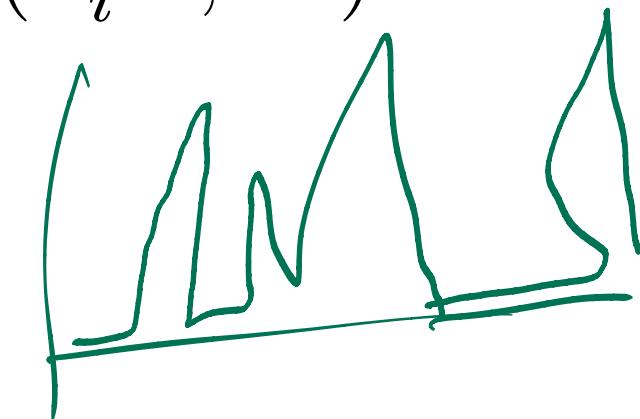
$$P(\hat{y} | x)$$

$$y = n \cdot x + b$$

- Now assume that $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$, then $y_i \sim \mathcal{N}(\mathbf{x}_i^T \boldsymbol{\theta}, \sigma^2)$
- We can write the conditional distribution as :



$$P(y_i | \mathbf{x}_i) \sim \mathcal{N}(0, \sigma^2)$$



Probabilistic Interpretation of LR

- Let's assume that all data points in the dataset are i.i.d. (independently identically distributed). Then we have:

$$\text{IP}(\mathcal{D}) = \prod_{i=1}^N \text{IP}(\underline{x_i}, y_i)$$

- Using Bayes Theorem we can write:

$$\prod_{i=1}^N \text{IP}(\underline{x_i}, y_i) = \prod_{i=1}^N \text{IP}(x_i) \text{IP}(y_i | x_i)$$

$p(x_i)$
 y_i
 x_i

P(the teacher drinks water) =

$P(\text{the})P(\text{teacher})P(\text{drinks} \mid \text{the teacher})$.

$P(\text{water} \mid \text{the teacher})$.

Maximum Likelihood Estimator

- In simple words, given the Dataset we want to find the values of the unknown parameters which maximize the probability of the Dataset.
- Using the definition of the conditional distribution we have

$$\text{P}(y_i | \mathbf{x}_i) = \frac{1}{\sigma \sqrt{2\pi}} \exp \left(- \underbrace{(y_i - \mathbf{x}_i^T \boldsymbol{\theta})}_{\text{Error term}} \right)$$

- Using the definition we get

$$\prod_{i=1}^N \text{P}(\mathbf{x}_i, y_i) = \prod_{i=1}^N \text{P}(\mathbf{x}_i) \prod_{i=1}^N \frac{1}{\sigma \sqrt{2\pi}} \exp \left(- \underbrace{(y_i - \mathbf{x}_i^T \boldsymbol{\theta})}_{\text{Error term}} \right)$$

Maximum Likelihood Estimator

- Let's try to maximaize:

$$\prod_{i=1}^N \text{P}(\mathbf{x}_i, y_i) = \prod_{i=1}^N \text{P}(\mathbf{x}_i) \prod_{i=1}^N \frac{1}{\sigma \sqrt{2\pi}} \exp(-(y_i - \mathbf{x}_i^T \boldsymbol{\theta}))$$

- Note that

$$\arg \max_{\boldsymbol{\theta}} \prod_{i=1}^N \text{P}(\mathbf{x}_i, y_i) \stackrel{\text{argmle}}{=} \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^N \exp(-(y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2)$$

Maximum Likelihood Estimator

- Furthermore, since the right hand side of the above equation is monotonic in \theta the arg max will not change if we take log of the expression

$$\arg \max_{\theta} \prod_{i=1}^N \exp(-(y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2) = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N (- (y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2)$$

argmin
 in interpretation
 of maximum
 likelihood

- Notice that the right hand side is minimising the MSE.
- Hence solution of minimizing the MSE is equivalent to Maximum Likelihood Estimator for linear regression

A Slight Detour: A Look at Optimization Tools

Introduction

- **Optimizers** are algorithms that adjust the weights of the model to minimize the loss function during training.
- They are one of the main components of Deep Learning models.
- We will go through several types of optimizers in this section.

Direction of maximum increase and decrease for a function

- Gradient direction is the direction of maximum increase for a function
- Negative gradient is the direction of maximum decrease for a function

$$\hat{w} = w - \eta w$$

```
# Vanilla Gradient Descent

while True:
    weights_grad = evaluate_gradient(loss_fun, data, weights)
    weights += - step_size * weights_grad # perform parameter update
```

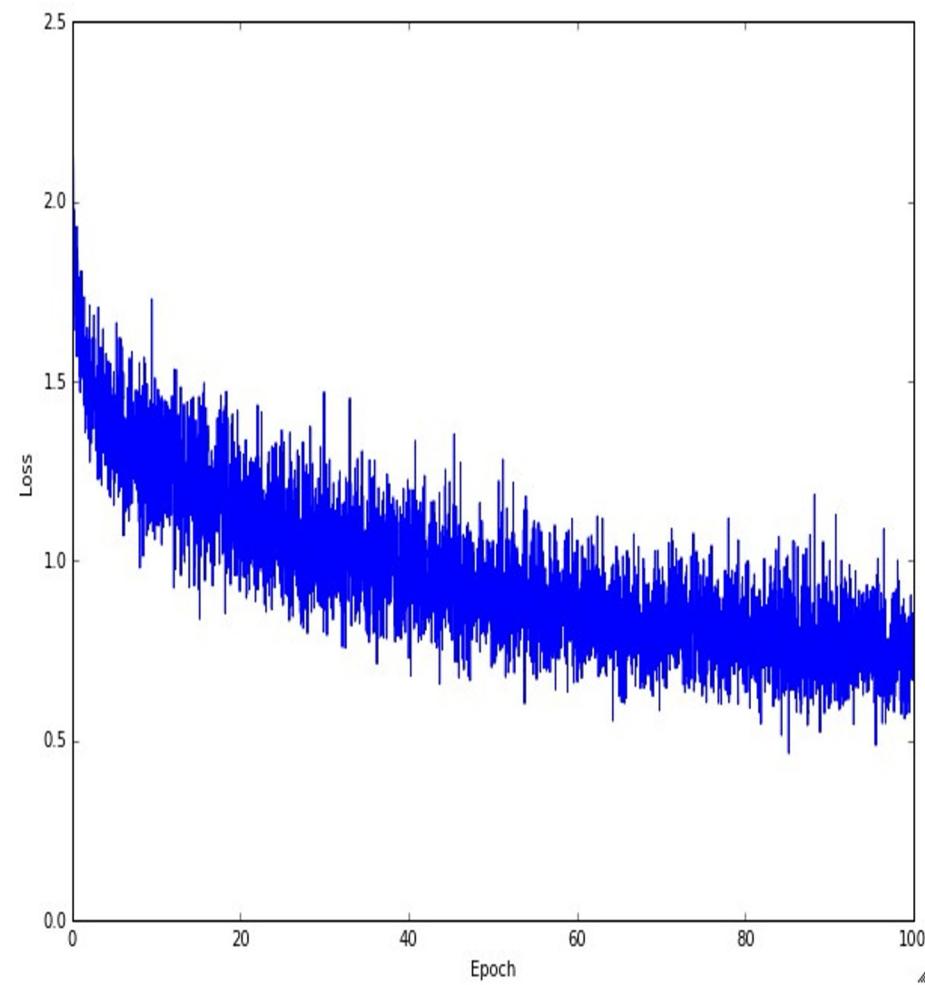
- **Mini batch Gradient Descent**
- only use a small portion of the training set to compute the gradient.

```
# Vanilla Minibatch Gradient Descent

while True:
    data_batch = sample_training_data(data, 256) # sample 256 examples
    weights_grad = evaluate_gradient(loss_fun, data_batch, weights)
    weights += - step_size * weights_grad # perform parameter update
```

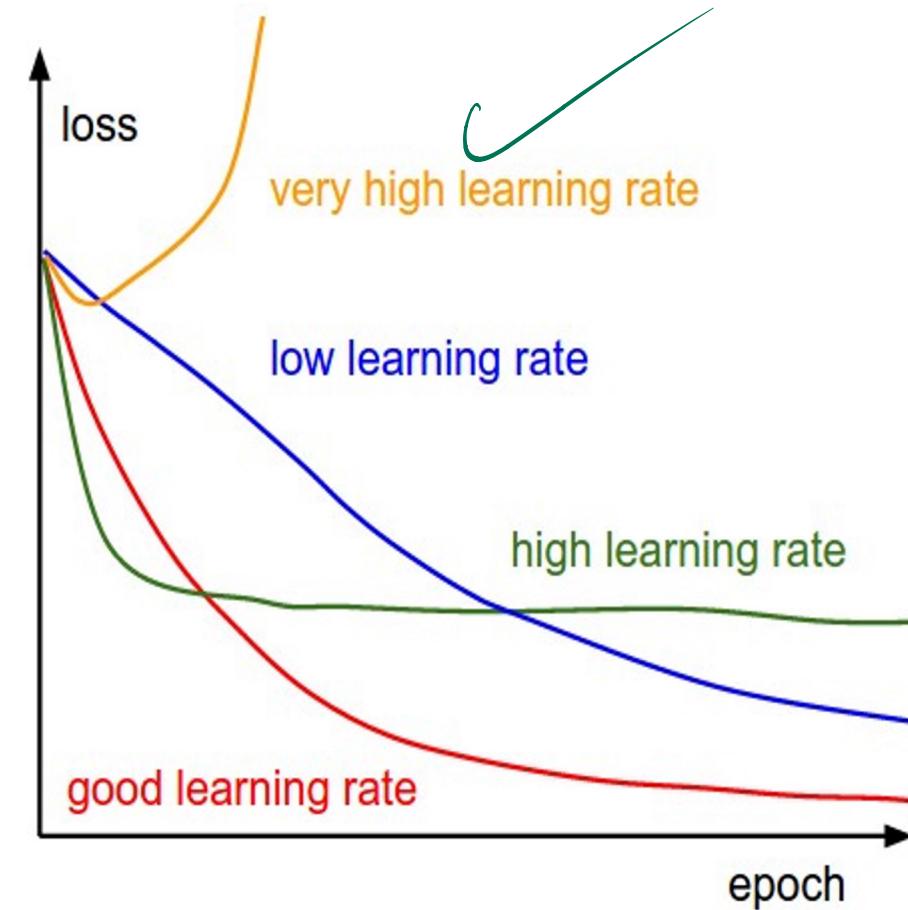
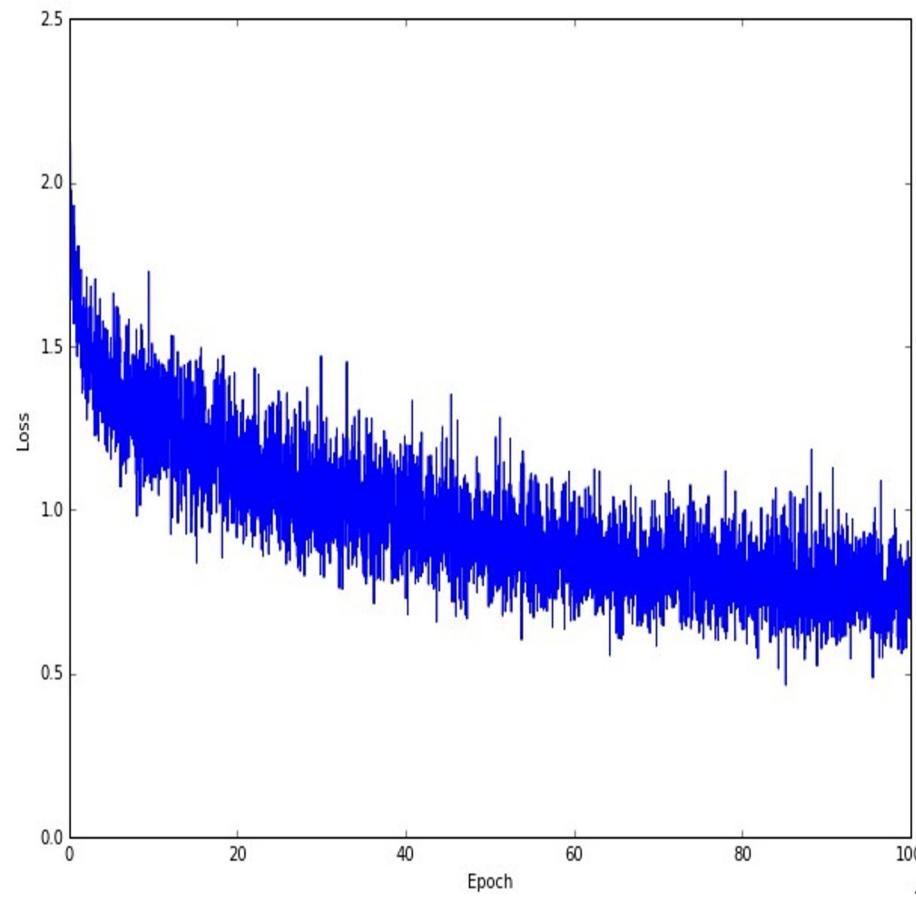
Common mini-batch sizes are 32/64/128 examples
e.g. Krizhevsky ILSVRC ConvNet used 256 examples

2^5



Example of optimization progress
while training a neural network.

(Loss over mini-batches goes
down over time.)



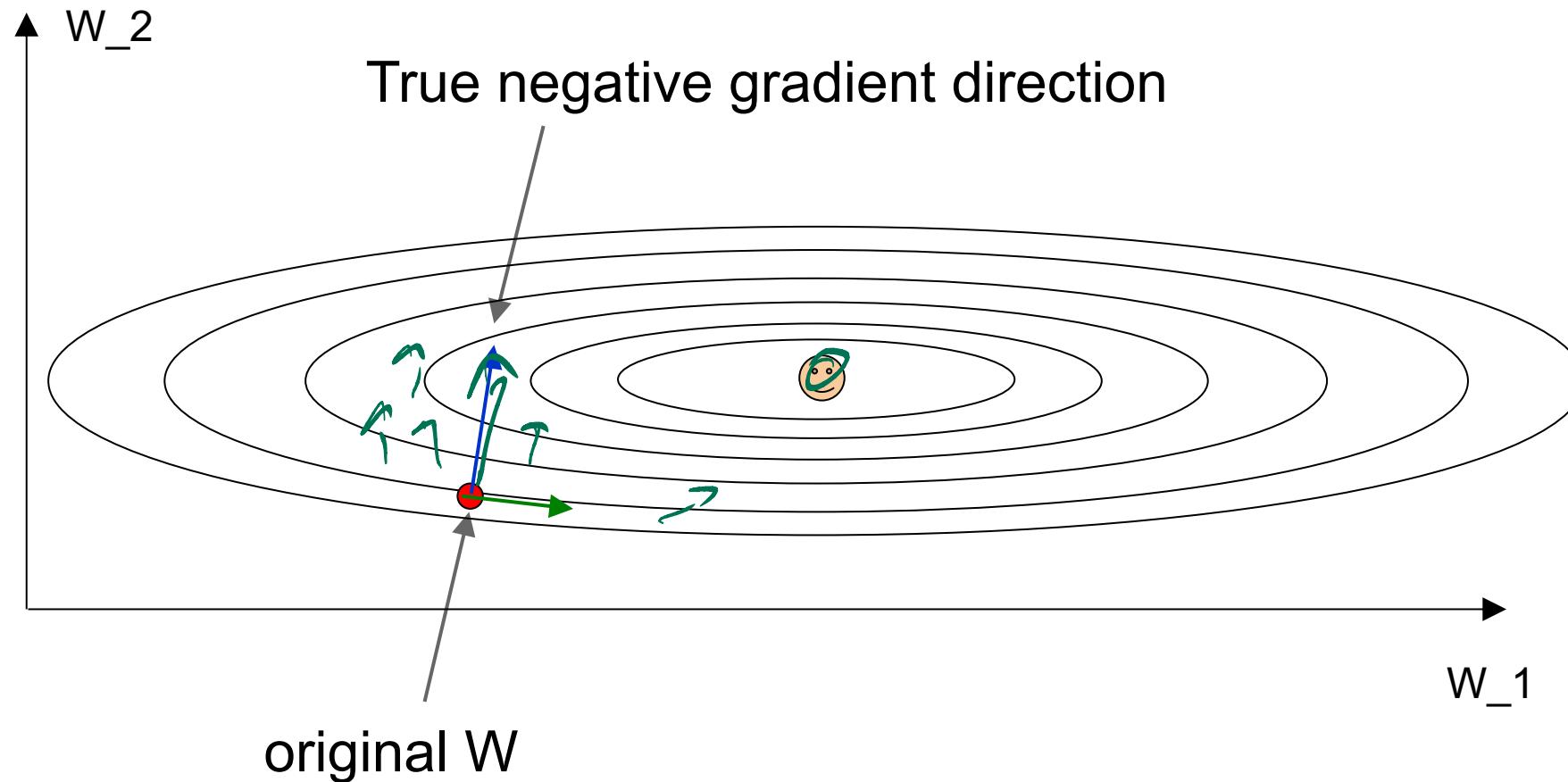
- only use a small portion of the training set to compute the gradient.

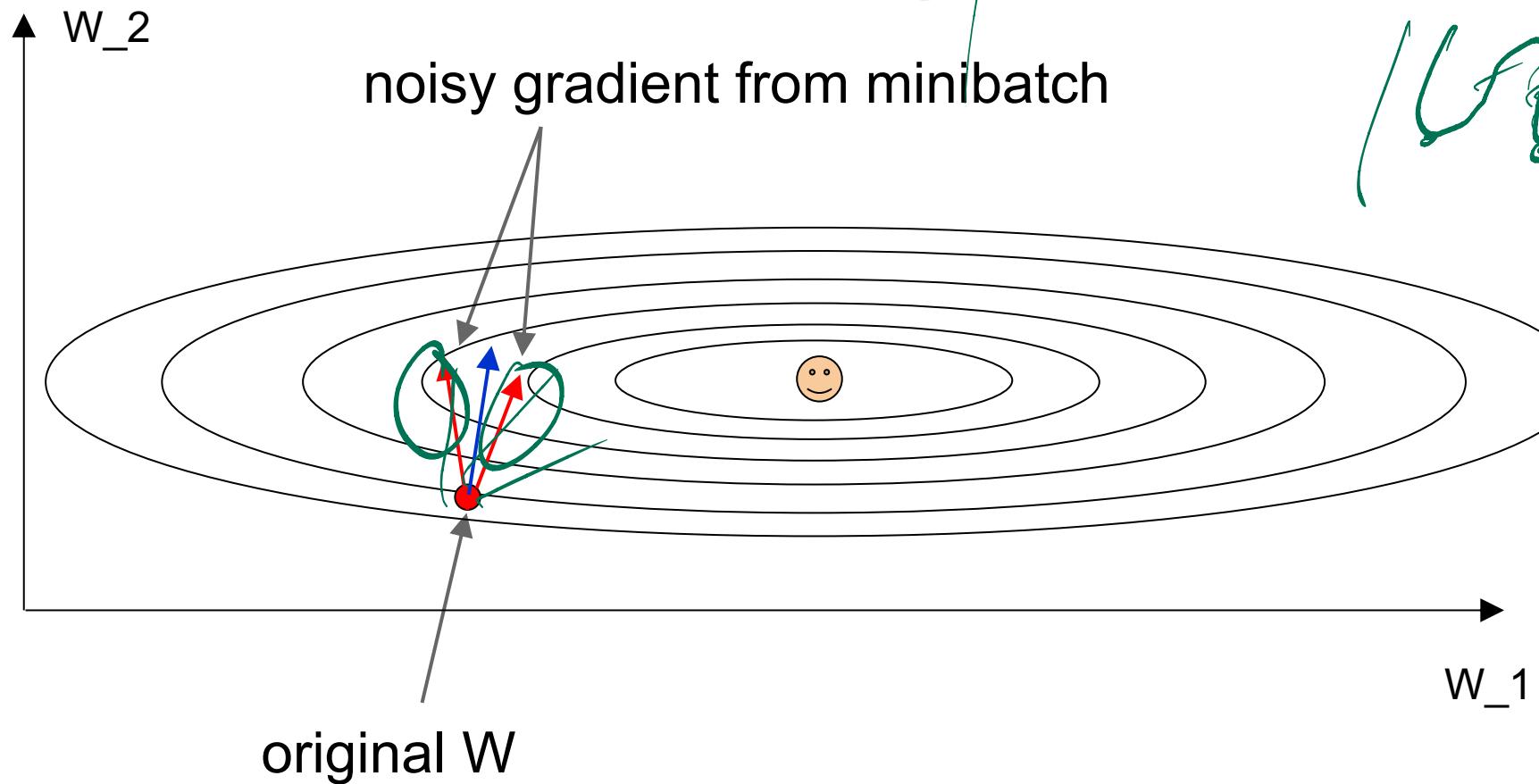
```
# Vanilla Minibatch Gradient Descent

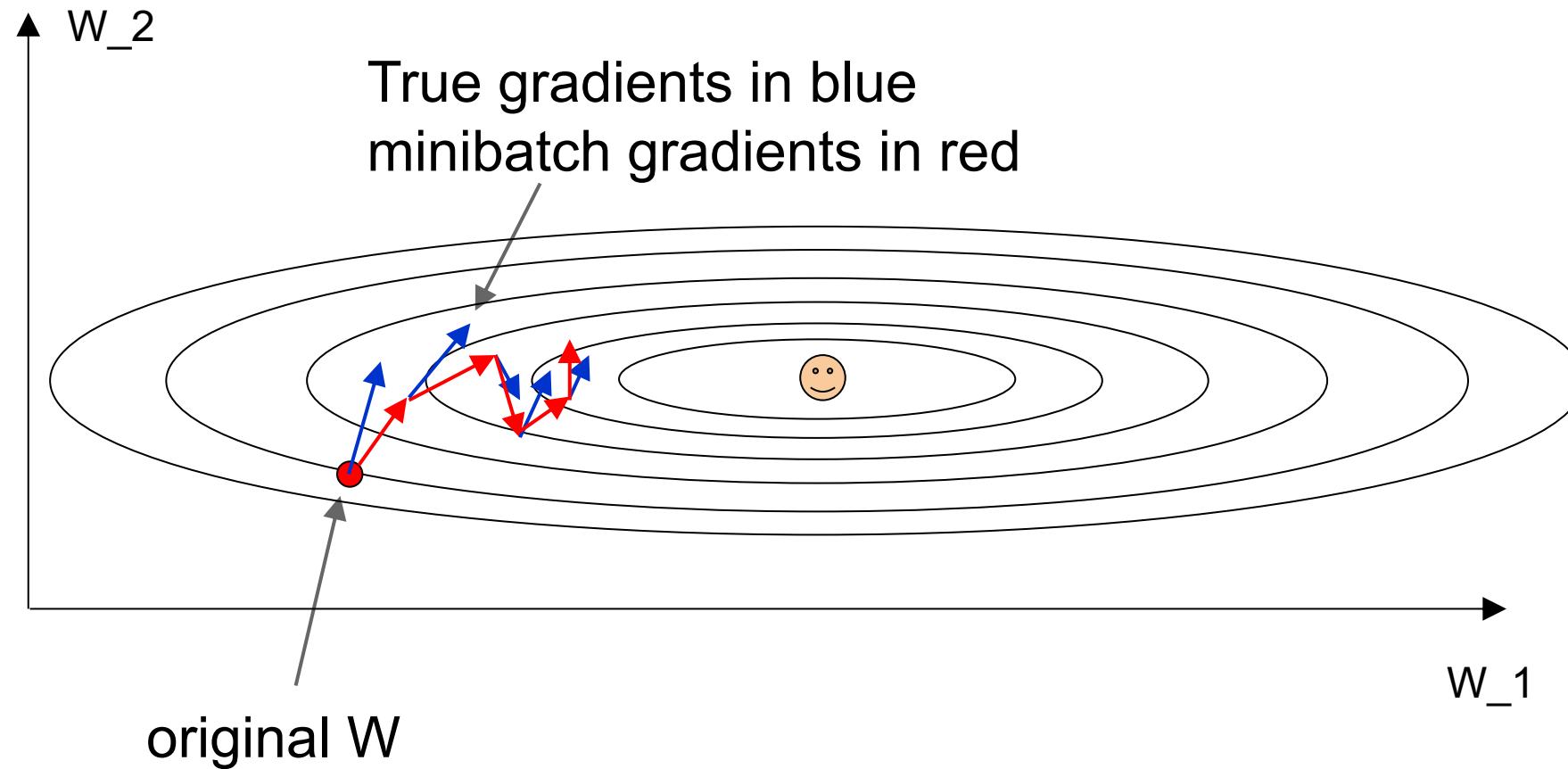
while True:
    data_batch = sample_training_data(data, 256) # sample 256 examples
    weights_grad = evaluate_gradient(loss_fun, data_batch, weights)
    weights += - step_size * weights_grad # perform parameter update
```

Common mini-batch sizes are 32/64/128 examples
e.g. Krizhevsky ILSVRC ConvNet used 256 examples

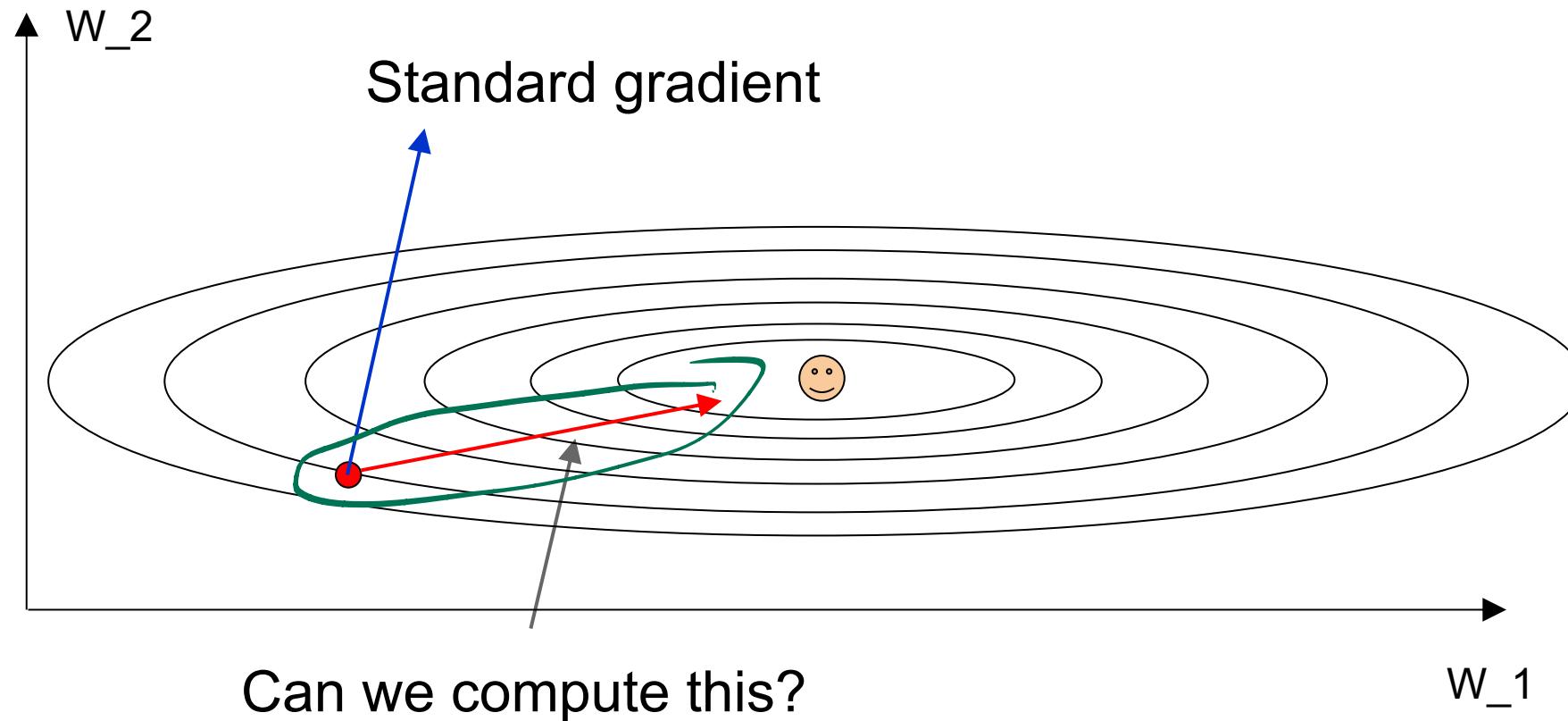
we will look at more fancy update formulas (momentum, Adagrad, RMSProp, Adam, ...)

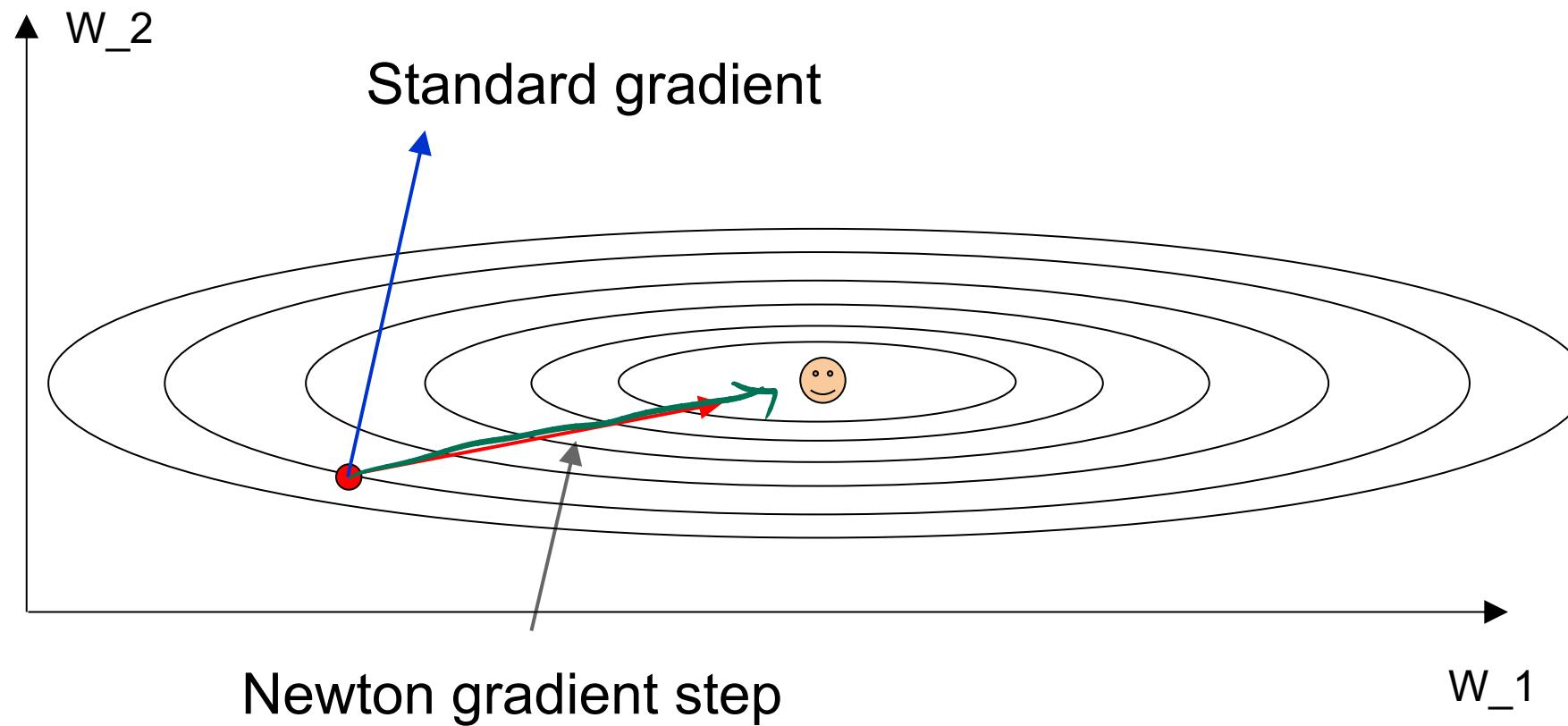






Gradients are noisy but still make good progress on average





higher

The Newton update is:

$$x_{n+1} = x_n - H_f(x_n)^{-1} \nabla f(x_n)$$

Where H_f is the Hessian matrix of second derivatives of f .

Converges very fast, but rarely used in DL. Why?

The Newton update is:

$$x_{n+1} = x_n - H_f(x_n)^{-1} \nabla f(x_n)$$

Converges very fast, but rarely used in DL. Why?

Too expensive: if x_n has dimension M, the Hessian $H_f(x_n)$ has dimension M^2 and takes $O(M^3)$ time to invert.

The Newton update is:

$$x_{n+1} = x_n - H_f(x_n)^{-1} \nabla f(x_n)$$

Converges very fast, but rarely used in DL. Why?

Too expensive: if x_n has dimension M, the Hessian $H_f(x_n)$ has dimension M^2 and takes $O(M^3)$ time to invert.

Too unstable: it involves a high-dimensional matrix inverse, which has poor numerical stability. The Hessian may even be singular.

Momentum Optimization

```
# Gradient descent update
x += - learning_rate * dx
```

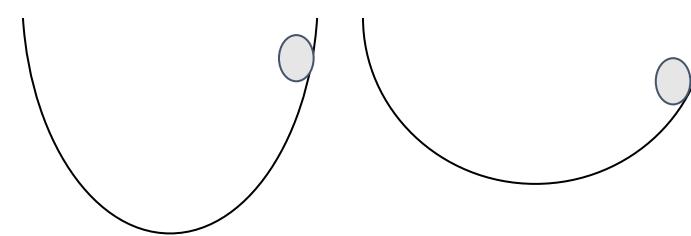


$\frac{1}{2}$

```
# Momentum update
v = mu * v - learning_rate * dx # integrate velocity
x += v # integrate position
```

- Physical interpretation as ball rolling down the loss function + friction (mu coefficient).
- mu = usually ~0.5, 0.9, or 0.99 (Sometimes annealed over time, e.g. from 0.5 -> 0.99) for adaptation

```
# Gradient descent update  
x += - learning_rate * dx
```



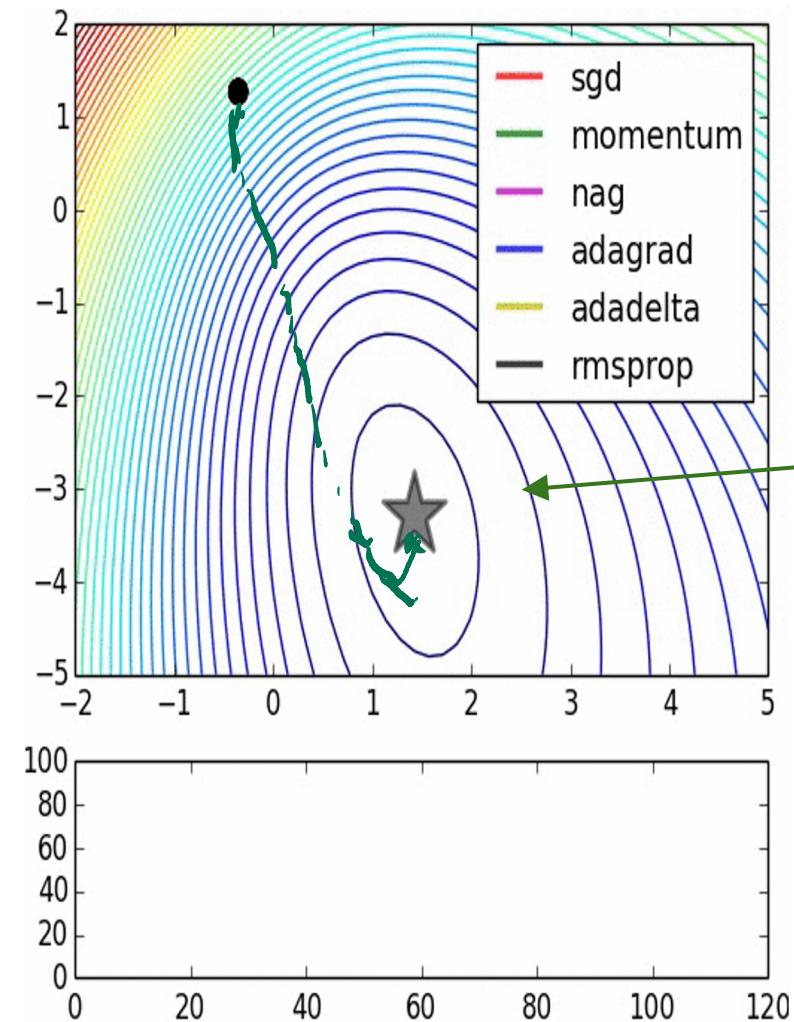
```
# Momentum update  
v = mu * v - learning_rate * dx # integrate velocity  
x += v # integrate position
```

- Allows a velocity to “build up” along shallow directions
- Velocity becomes damped in steep direction due to quickly changing sign



Smoother trajectories in optimization, avoiding zig-zagging behavior.
Faster convergence towards the global minimum.

SGD vs Momentum



notice momentum
overshooting the target,
but overall getting to the
minimum much faster
than vanilla SGD.

```
# Adagrad update
cache += dx**2
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```

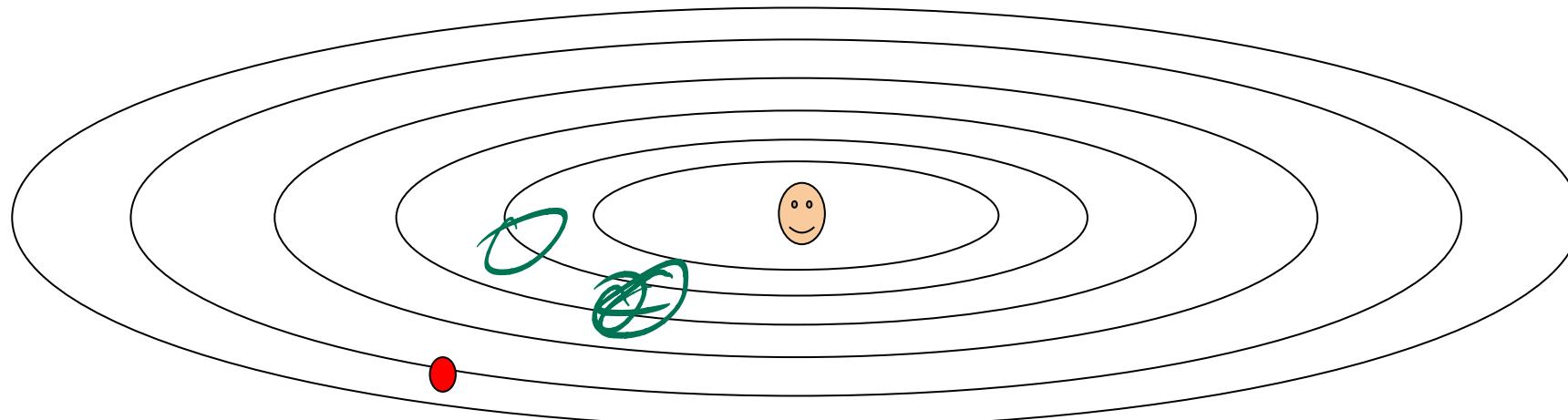
gradient clipping

Added element-wise scaling of the gradient based on the historical sum of squares in each dimension

- The element-wise scaling ensures that parameters corresponding to dimensions with larger cumulative gradients (cache) are updated less aggressively.
- Conversely, parameters with smaller cumulative gradients are updated more, balancing their influence.

making larger adjustments in directions where gradients are smaller (flatter regions) and smaller adjustments in steeper directions.

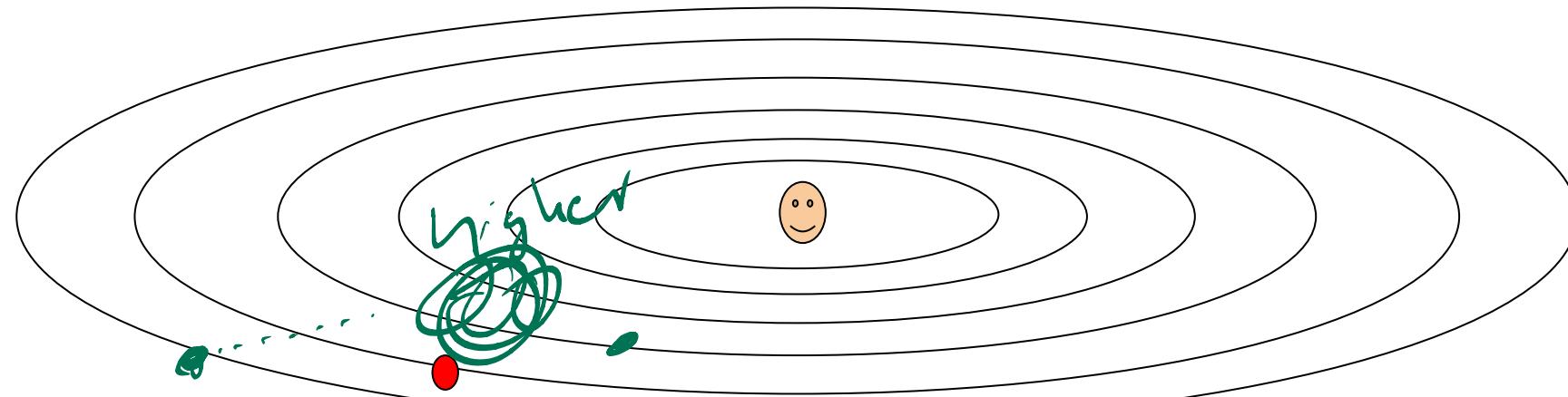
```
# Adagrad update
cache += dx**2
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```



Q: What happens with AdaGrad?

The optimizer takes smaller steps vertically (steep direction) and larger steps horizontally (flat direction).

```
# Adagrad update
cache += dx**2
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```



Q2: What happens to the step size over long time?

Extremely small steps as the optimizer nears convergence, effectively "stalling" progress.

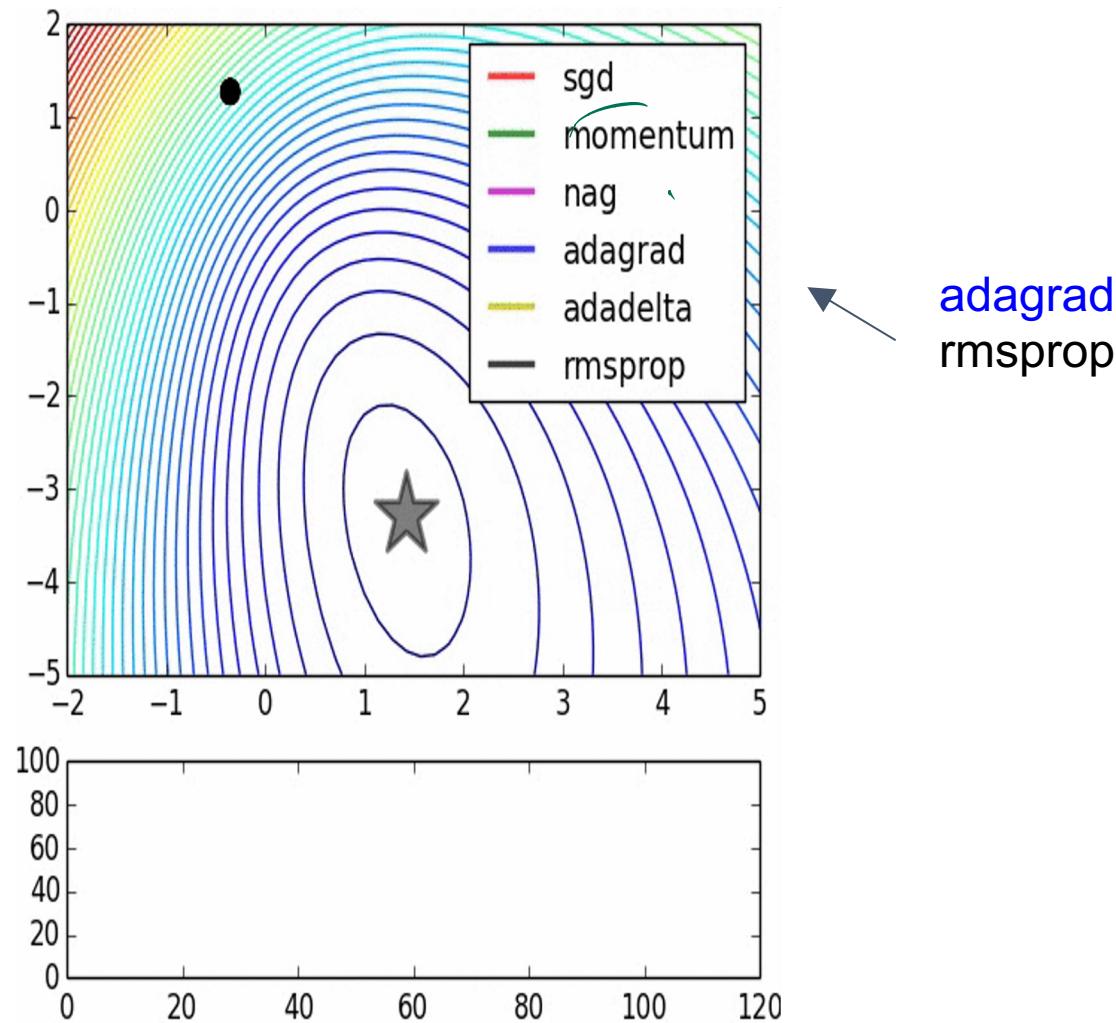
RMSProp

Exponential decay factor for cache, which prevents it from growing indefinitely.

```
# Adagrad update
cache += dx**2
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```



```
# RMSProp
cache = decay_rate * cache + (1 - decay_rate) * dx**2
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```



Adam optimizer

smooths the gradients over time by keeping a moving average of the gradient.

```
# Adam
m = beta1*m + (1-beta1)*dx # update first moment
v = beta2*v + (1-beta2)*(dx**2) # update second moment
x += - learning_rate * m / (np.sqrt(v) + 1e-7)
```

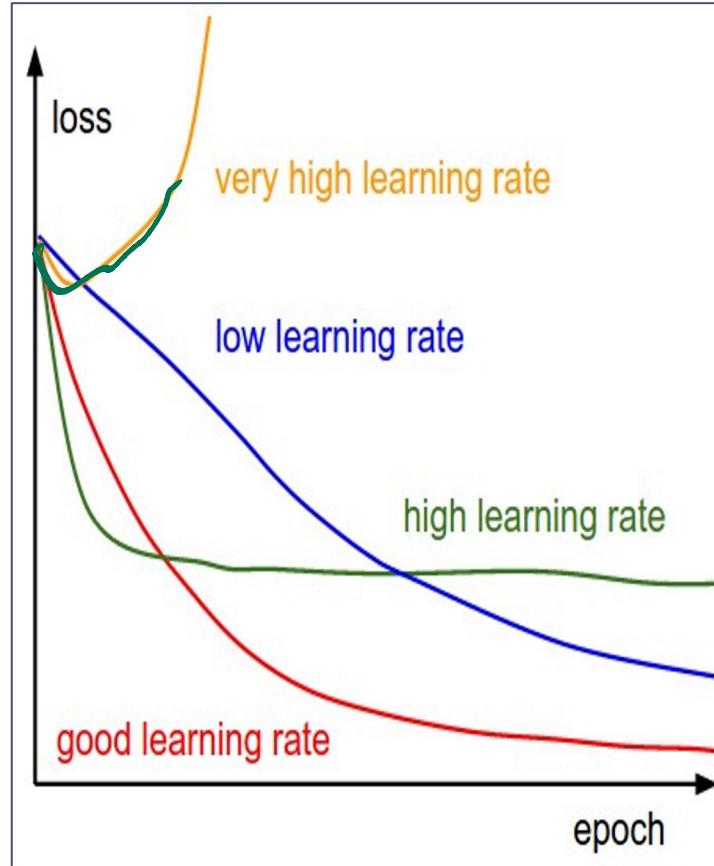
momentum

RMSProp-like

normalizing gradients by scaling them with the square root of their historical squared values.

Looks a bit like RMSProp with momentum

for



Q: Which one of these learning rates is best to use?

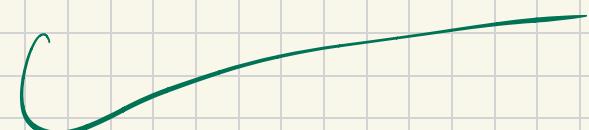
Small

model.train()

for i in range(100):

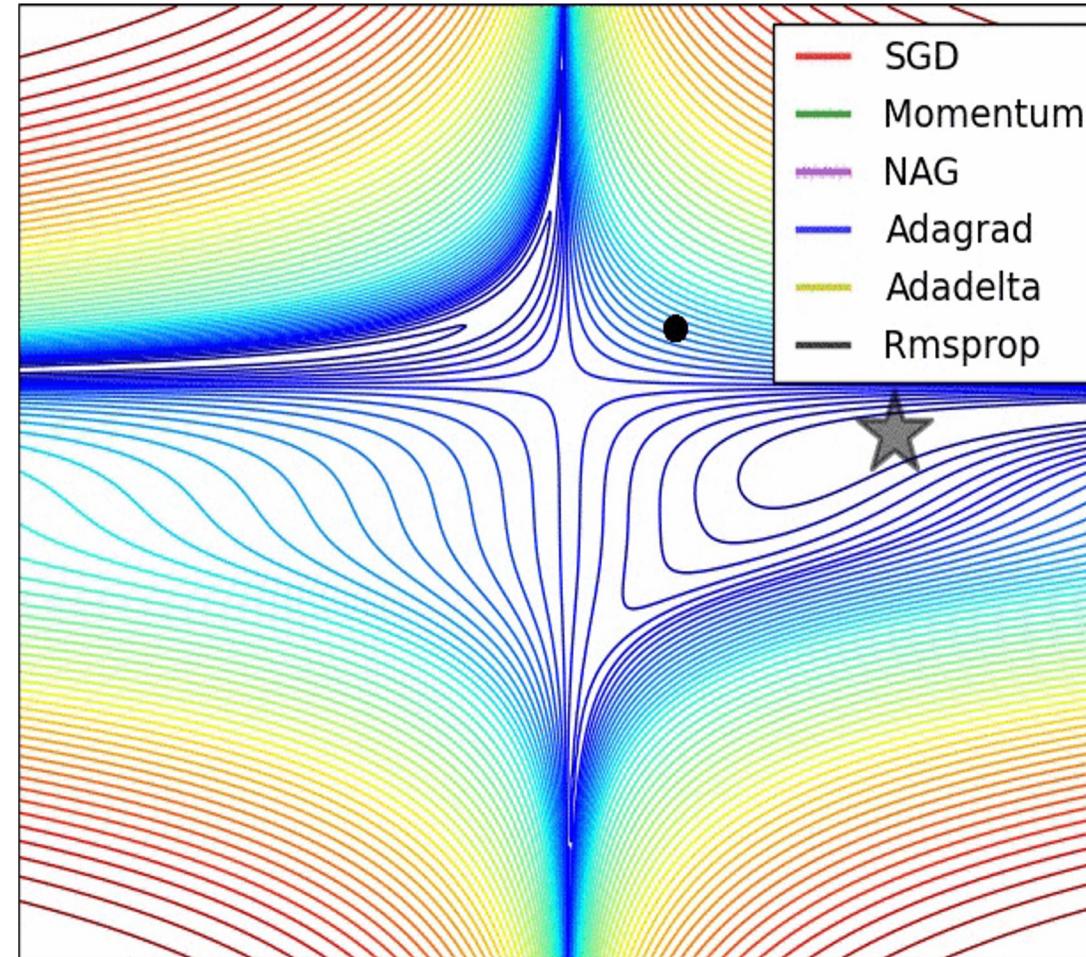
out = model(i)

$$\text{loss} = \text{CE}(out, y)$$



→ loss + if loss > 1:

If epoch -
adjust Q - .()



(image credits to Alec Radford)

Summary

