

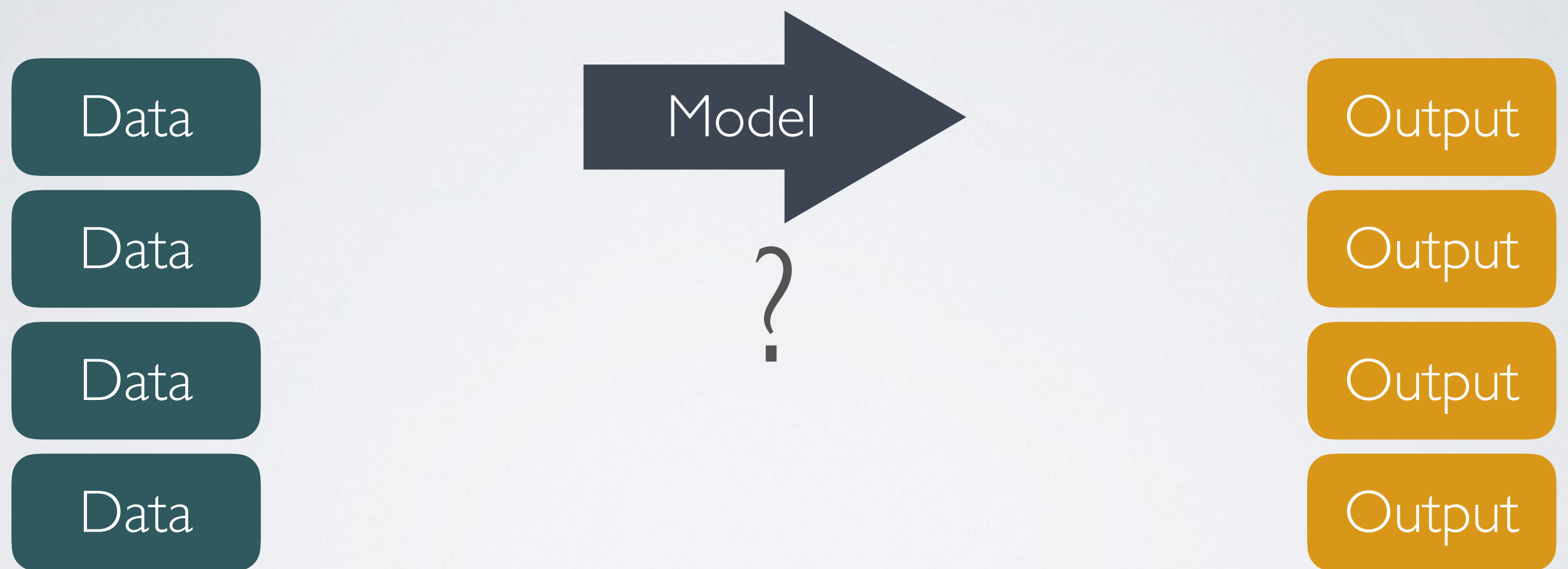
STOCHASTIC GRADIENT DESCENT

Philipp Krähenbühl

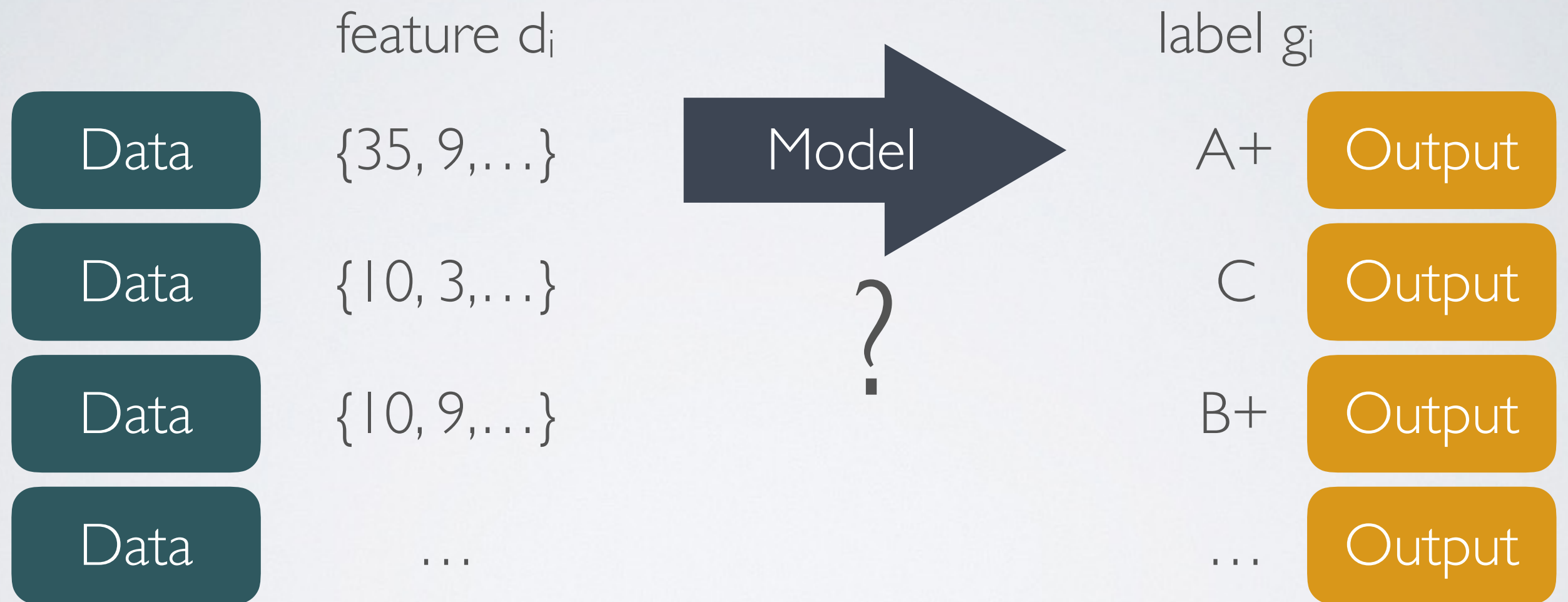
ADMINISTRATIVE

- List of papers due **Thursday 6am**
- Custer access (work in progress)
- Caffe tutorial for cs381 v **Sep 12** 5pm-7pm
 - more details to follow

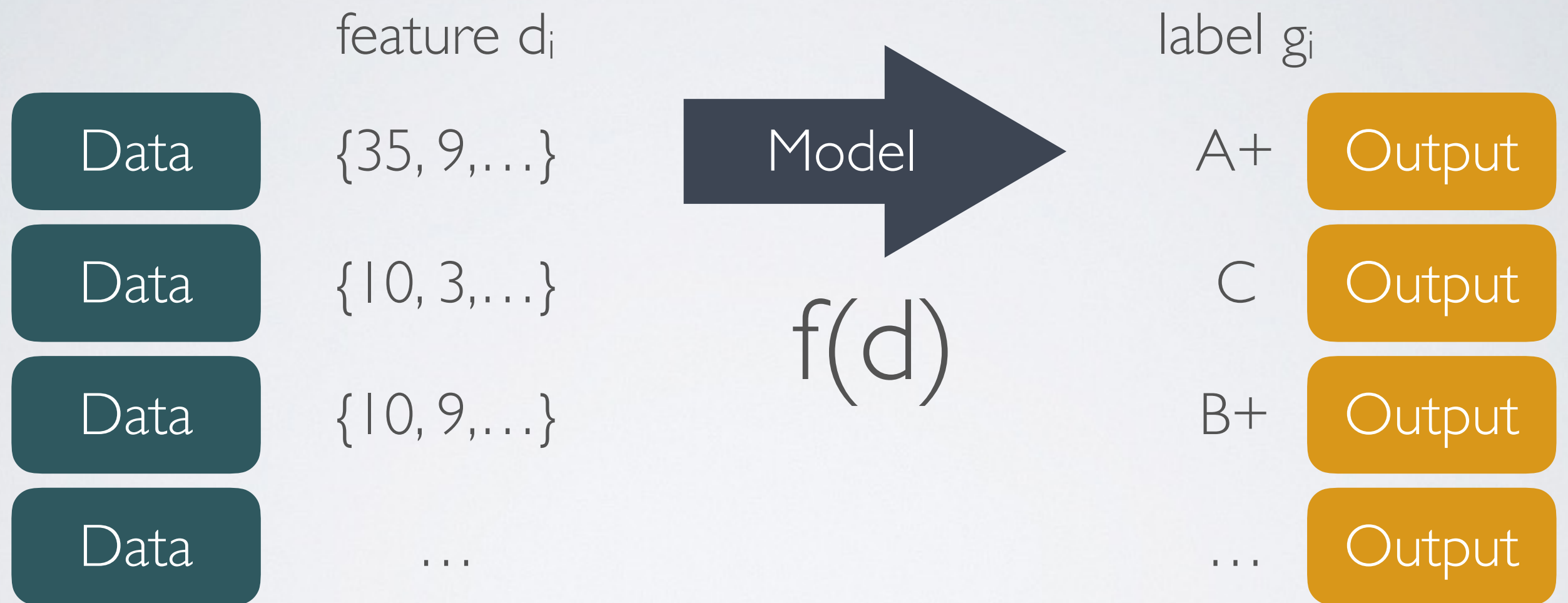
RECAP



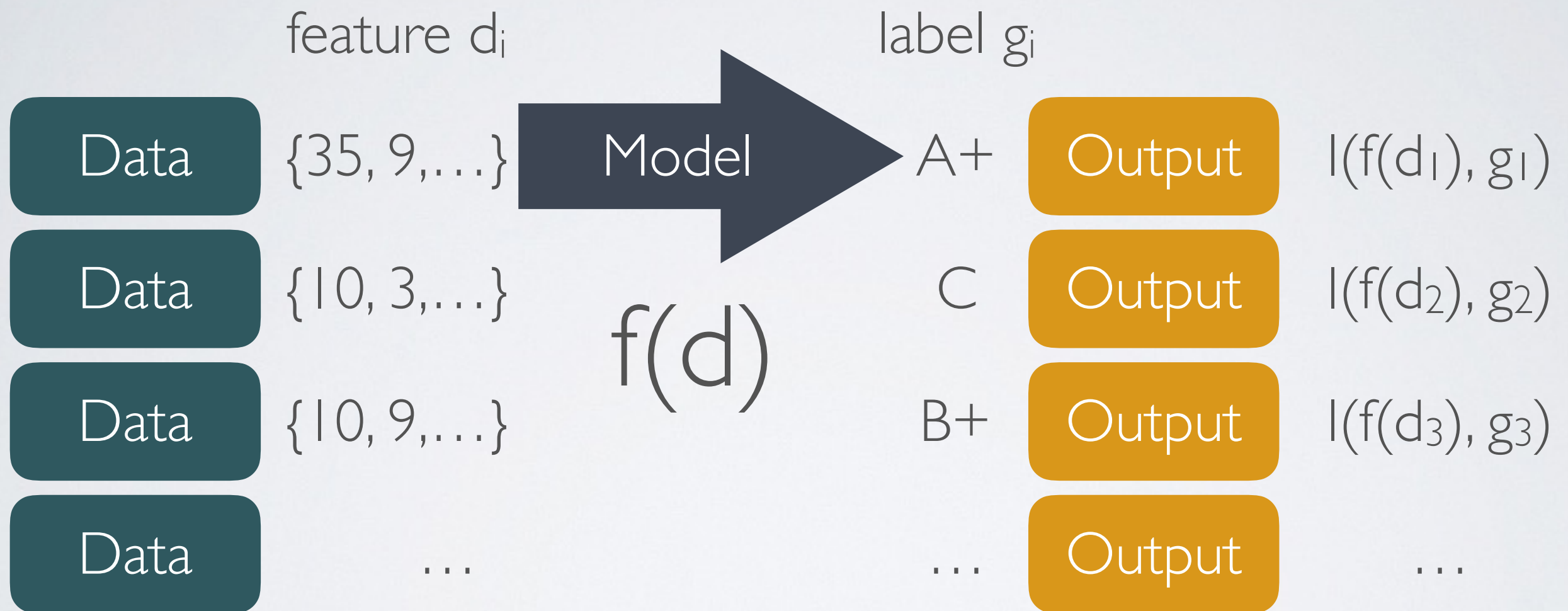
RECAP



RECAP



RECAP



RECAP

$$L = l(f(d_1), g_1) + l(f(d_2), g_2) + l(f(d_3), g_3) + \dots$$

RECAP

$$L = E_{(d,g) \sim D} [I(f(d), g)]$$

GRADIENT DESCENT

$$\nabla L = E_{(d,g) \sim D} [\nabla l(f(d), g)]$$

$$f = f - \alpha \nabla L$$

GRADIENT DESCENT

$l(f(d_1), g_1)$



$l(f(d_2), g_2)$

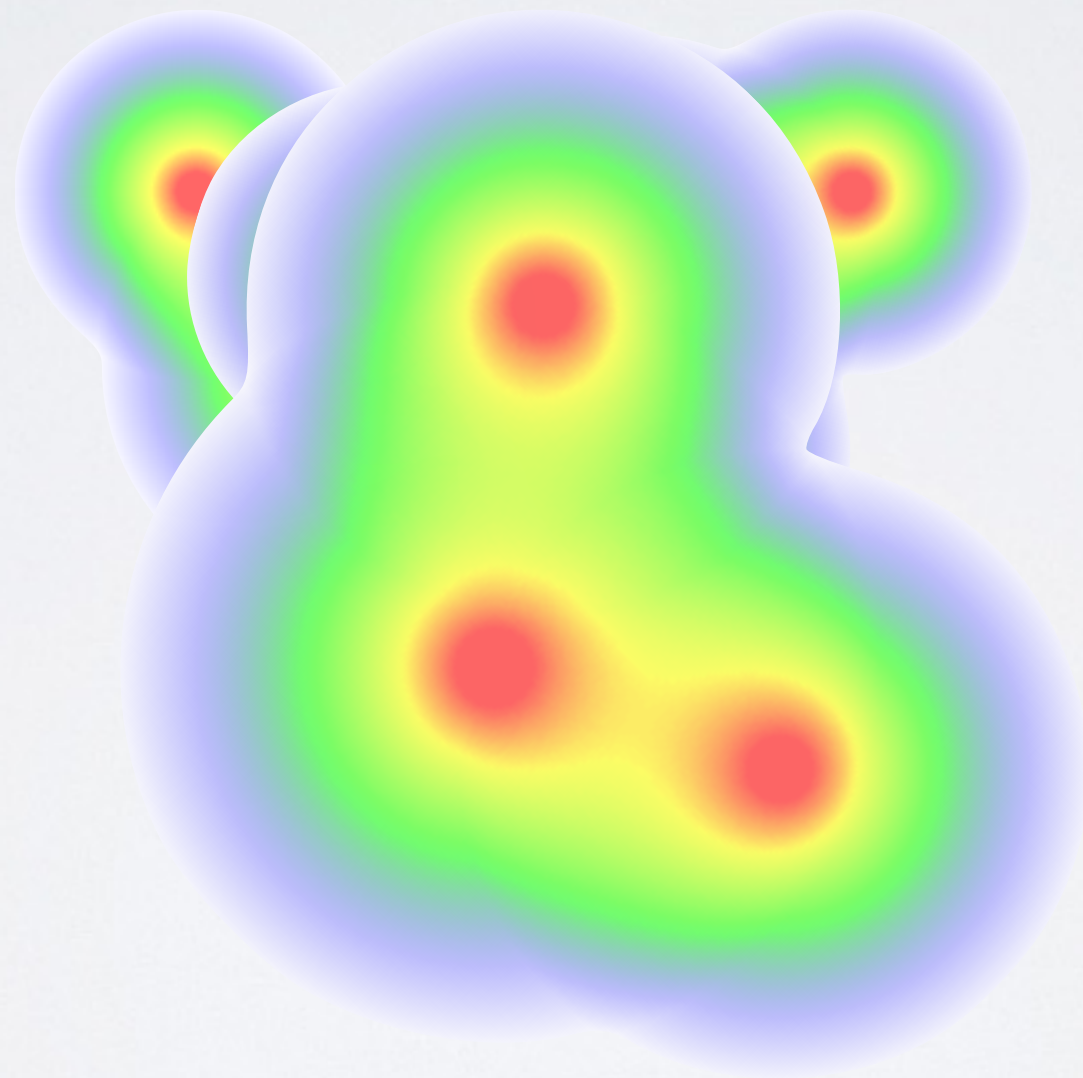


$l(f(d_3), g_3)$

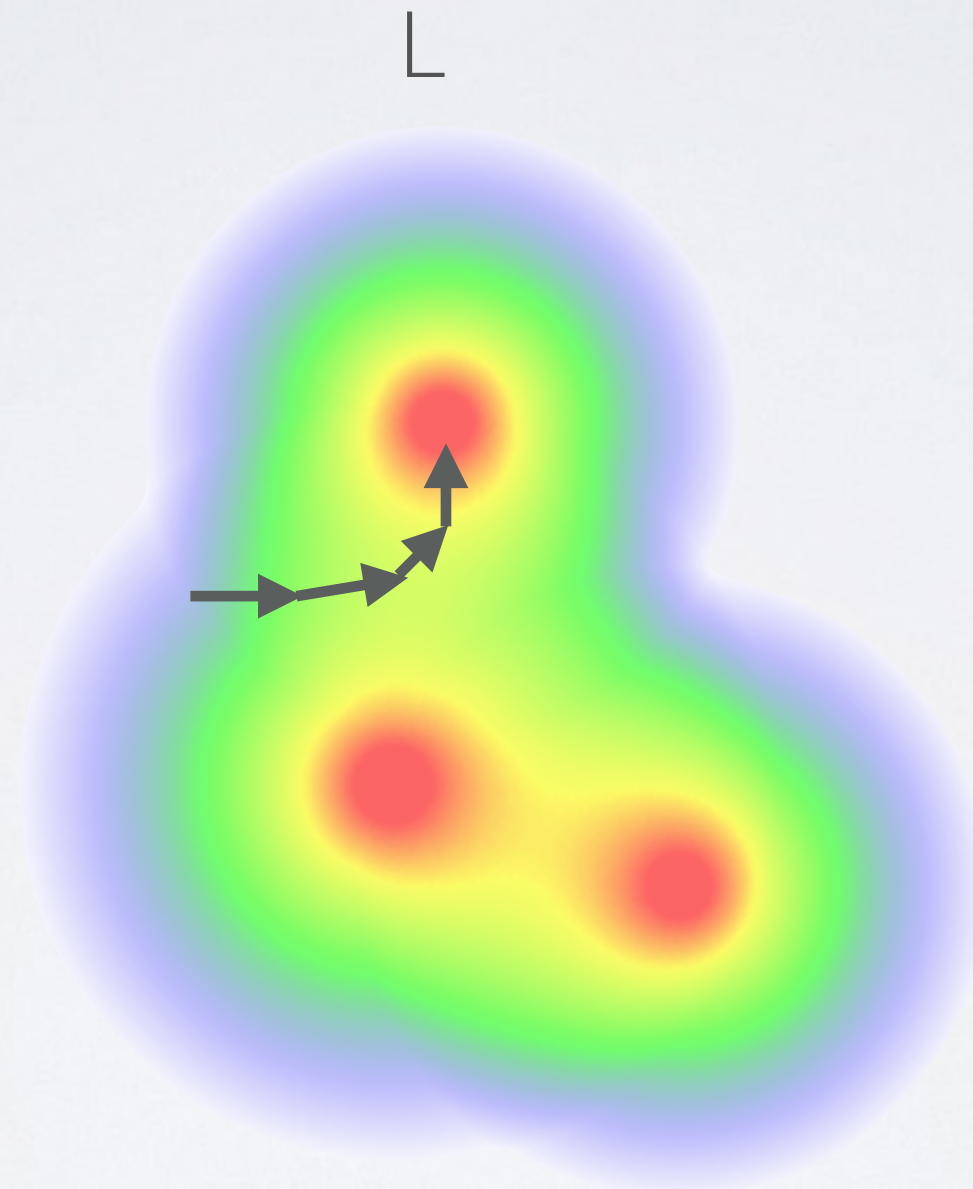


GRADIENT DESCENT

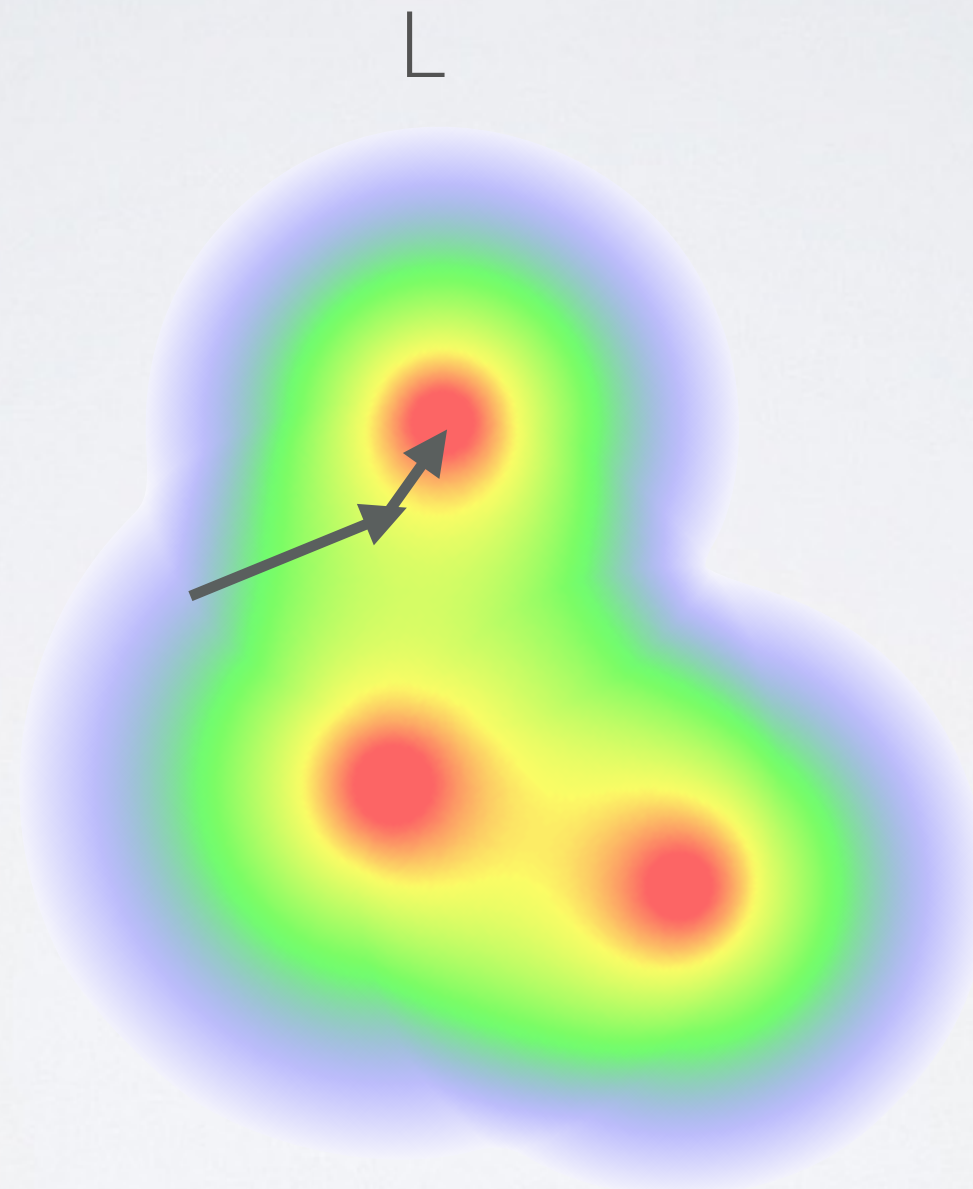
L



GRADIENT DESCENT

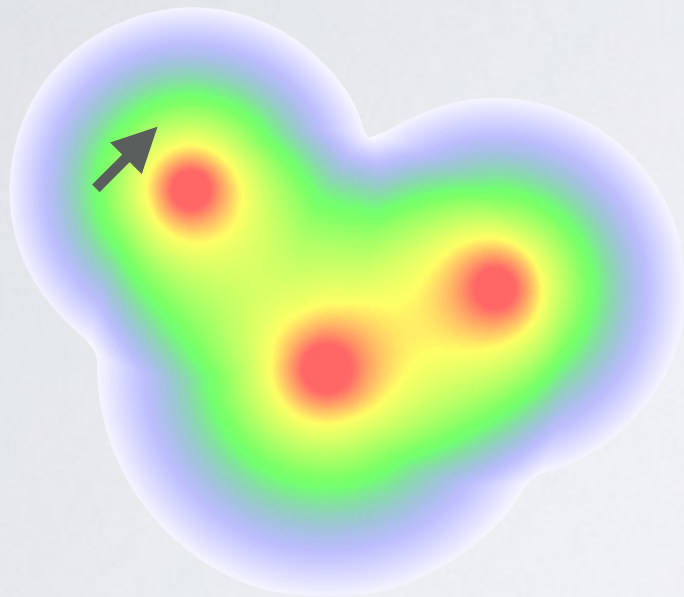


(2ND ORDER) GRADIENT DESCENT



STOCHASTIC GRADIENT DESCENT

$l(f(d_1), g_1)$



$l(f(d_2), g_2)$



$l(f(d_3), g_3)$



STOCHASTIC GRADIENT DESCENT

$$\mathbf{f} = \mathbf{f} - \gamma \nabla l(\mathbf{f}(\mathbf{d}), \mathbf{g}) \quad \forall (\mathbf{d}, \mathbf{g}) \sim D$$

SGD ANALYSIS

$$\begin{aligned} E(\tilde{f}_n) - E(f^*) &= E(f_{\mathcal{F}}^*) - E(f^*) && \text{Approximation error } (\mathcal{E}_{\text{app}}) \\ &+ E(f_n) - E(f_{\mathcal{F}}^*) && \text{Estimation error } (\mathcal{E}_{\text{est}}) \\ &+ E(\tilde{f}_n) - E(f_n) && \text{Optimization error } (\mathcal{E}_{\text{opt}}) \end{aligned}$$

Problem:

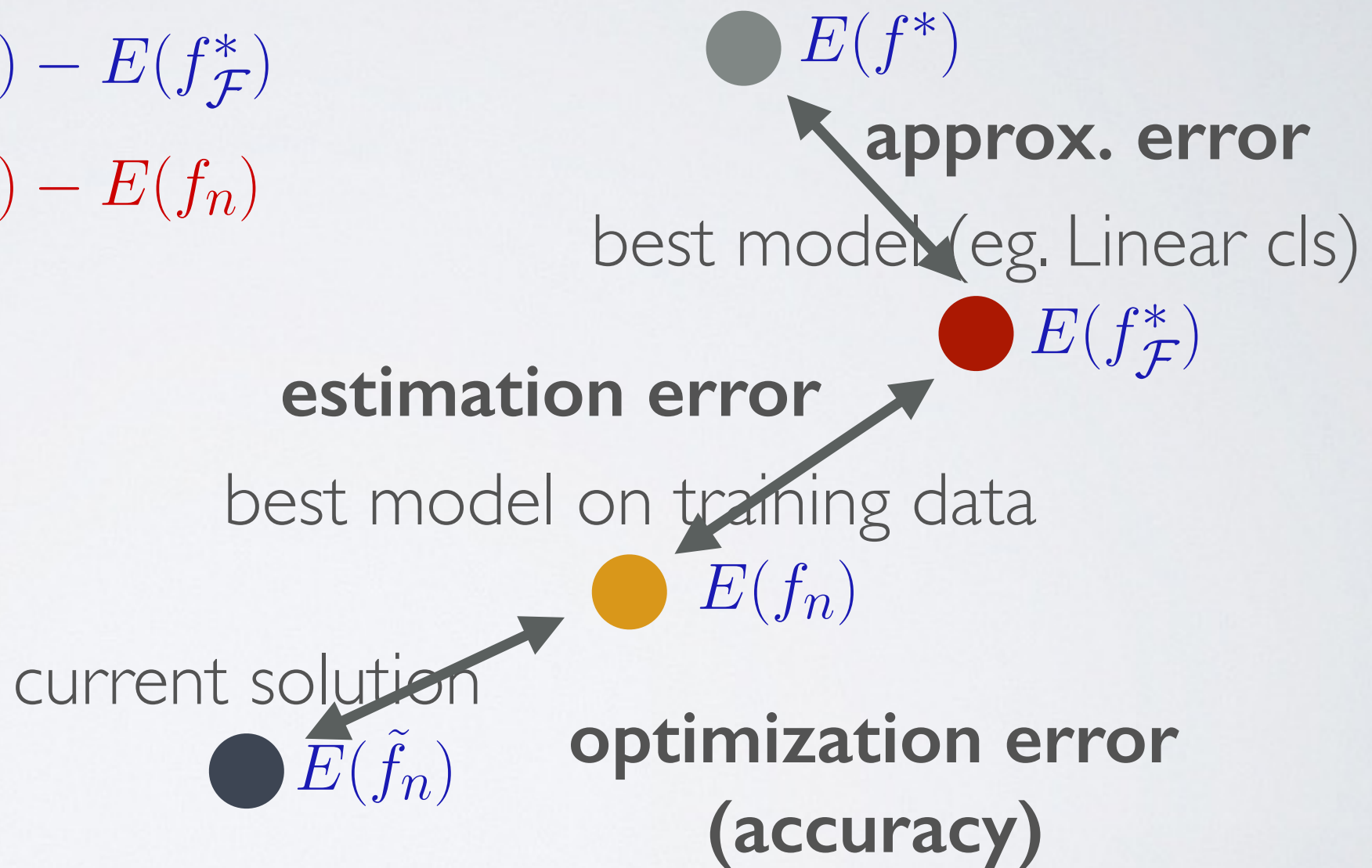
Choose \mathcal{F} , n , and ρ to make this as small as possible,

subject to budget constraints $\left\{ \begin{array}{l} \text{max number of examples } n \\ \text{max computing time } T \end{array} \right.$

SGD ANALYSIS

$$\begin{aligned} E(\tilde{f}_n) - E(f^*) &= E(f_{\mathcal{F}}^*) - E(f^*) \\ &+ E(f_n) - E(f_{\mathcal{F}}^*) \\ &+ E(\tilde{f}_n) - E(f_n) \end{aligned}$$

true underlying model



SMALL SCALE LEARNING

$$\begin{aligned} E(\tilde{f}_n) - E(f^*) &= E(f_{\mathcal{F}}^*) - E(f^*) \\ &+ E(f_n) - E(f_{\mathcal{F}}^*) \\ &+ E(\tilde{f}_n) - E(f_n) \end{aligned}$$

true underlying model

● $E(f^*)$

approx. error

best model (eg. Linear cls)

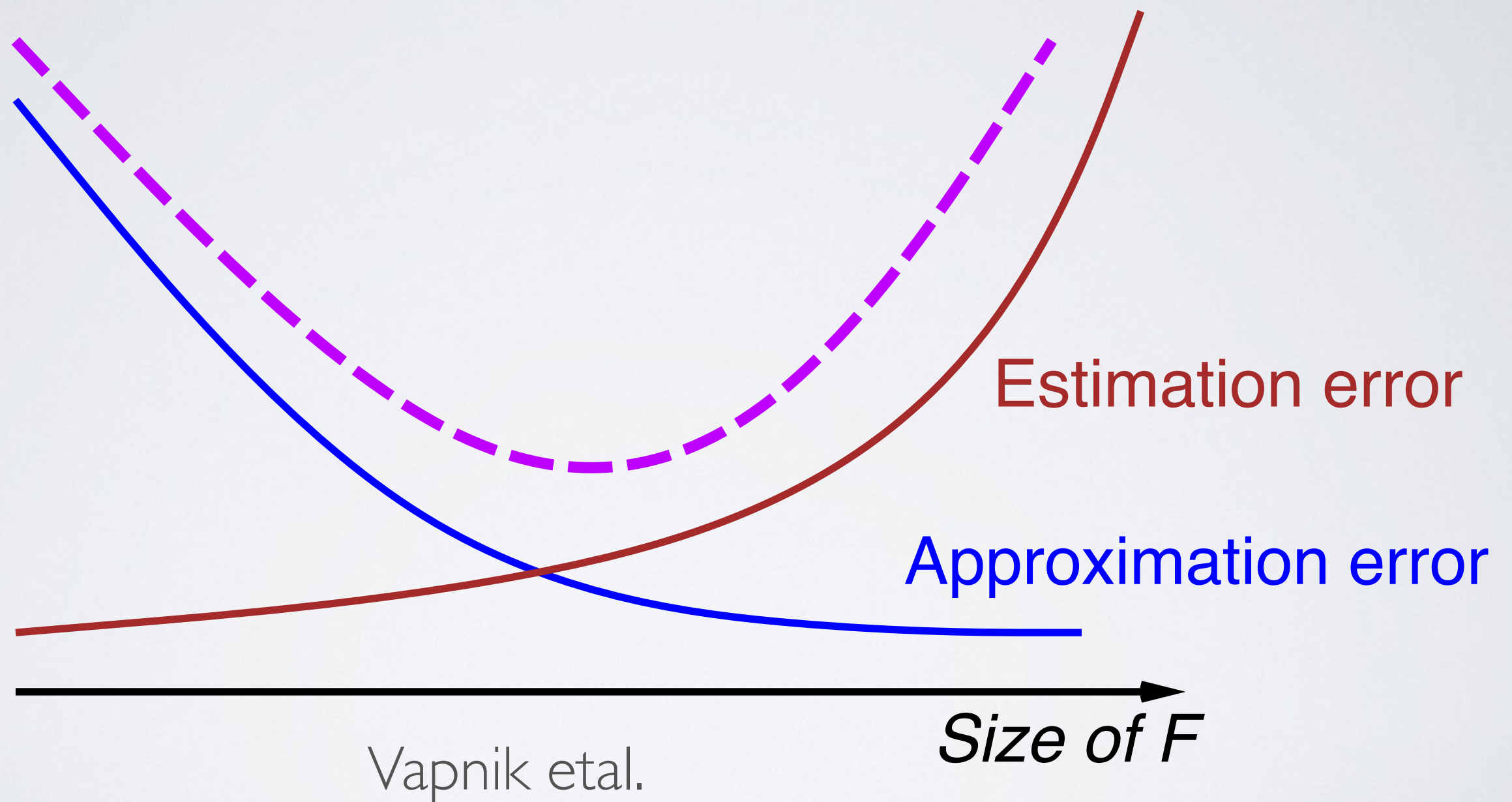
● $E(f_{\mathcal{F}}^*)$

estimation error

current solution = best model on training data

● $E(f_n) = E(\tilde{f}_n)$

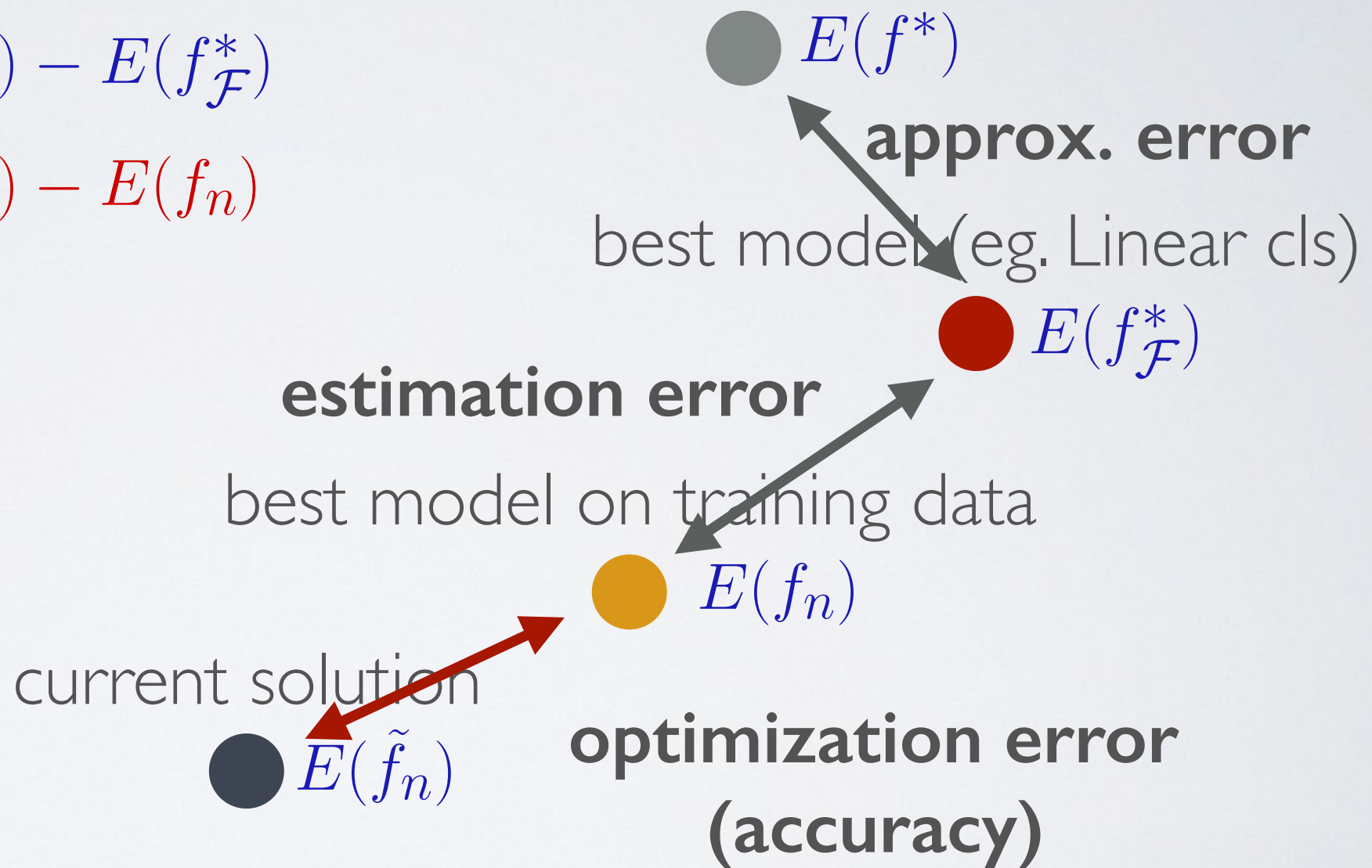
SMALL SCALE LEARNING



SGD ANALYSIS

$$\begin{aligned} E(\tilde{f}_n) - E(f^*) &= E(f_{\mathcal{F}}^*) - E(f^*) \\ &+ E(f_n) - E(f_{\mathcal{F}}^*) \\ &+ E(\tilde{f}_n) - E(f_n) \end{aligned}$$

true underlying model



SGD ANALYSIS

$$\begin{aligned} E(\tilde{f}_n) - E(f^*) &= E(f_{\mathcal{F}}^*) - E(f^*) && \text{Approximation error } (\mathcal{E}_{\text{app}}) \\ &+ E(f_n) - E(f_{\mathcal{F}}^*) && \text{Estimation error } (\mathcal{E}_{\text{est}}) \\ &+ E(\tilde{f}_n) - E(f_n) && \text{Optimization error } (\mathcal{E}_{\text{opt}}) \end{aligned}$$

Problem:

Choose \mathcal{F} , n , and ρ to make this as small as possible,

subject to budget constraints $\left\{ \begin{array}{l} \text{max number of examples } n \\ \text{max computing time } T \end{array} \right.$

STOCHASTIC GRADIENT DESCENT

	GD	2GD	SGD
Time per iteration :	n	n	1
Iterations to accuracy ρ :	$\log \frac{1}{\rho}$	$\log \log \frac{1}{\rho}$	$\frac{1}{\rho}$
Time to accuracy ρ :	$n \log \frac{1}{\rho}$	$n \log \log \frac{1}{\rho}$	$\frac{1}{\rho}$
Time to excess error ε :	$\frac{1}{\varepsilon^{1/\alpha}} \log^2 \frac{1}{\varepsilon}$	$\frac{1}{\varepsilon^{1/\alpha}} \log \frac{1}{\varepsilon} \log \log \frac{1}{\varepsilon}$	$\frac{1}{\varepsilon}$

STOCHASTIC GRADIENT DESCENT - EXAMPLES

Loss	Stochastic gradient algorithm
Adaline (Widrow and Hoff, 1960) $Q_{\text{adaline}} = \frac{1}{2} (y - w^\top \Phi(x))^2$ $\Phi(x) \in \mathbb{R}^d, y = \pm 1$	$w \leftarrow w + \gamma_t (y_t - w^\top \Phi(x_t)) \Phi(x_t)$
Perceptron (Rosenblatt, 1957) $Q_{\text{perceptron}} = \max\{0, -y w^\top \Phi(x)\}$ $\Phi(x) \in \mathbb{R}^d, y = \pm 1$	$w \leftarrow w + \gamma_t \begin{cases} y_t \Phi(x_t) & \text{if } y_t w^\top \Phi(x_t) \leq 0 \\ 0 & \text{otherwise} \end{cases}$
K-Means (MacQueen, 1967) $Q_{\text{kmeans}} = \min_k \frac{1}{2} (z - w_k)^2$ $z \in \mathbb{R}^d, w_1 \dots w_k \in \mathbb{R}^d$ $n_1 \dots n_k \in \mathbb{N}$, initially 0	$k^* = \arg \min_k (z_t - w_k)^2$ $n_{k^*} \leftarrow n_{k^*} + 1$ $w_{k^*} \leftarrow w_{k^*} + \frac{1}{n_{k^*}} (z_t - w_{k^*})$

STOCHASTIC GRADIENT DESCENT - EXAMPLES

SVM (Cortes and Vapnik, 1995)

$$Q_{\text{svm}} = \lambda w^2 + \max\{0, 1 - y w^\top \Phi(x)\} \quad w \leftarrow w - \gamma_t \begin{cases} \lambda w & \text{if } y_t w^\top \Phi(x_t) > 1, \\ \lambda w - y_t \Phi(x_t) & \text{otherwise.} \end{cases}$$
$$\Phi(x) \in \mathbb{R}^d, \quad y = \pm 1, \quad \lambda > 0$$

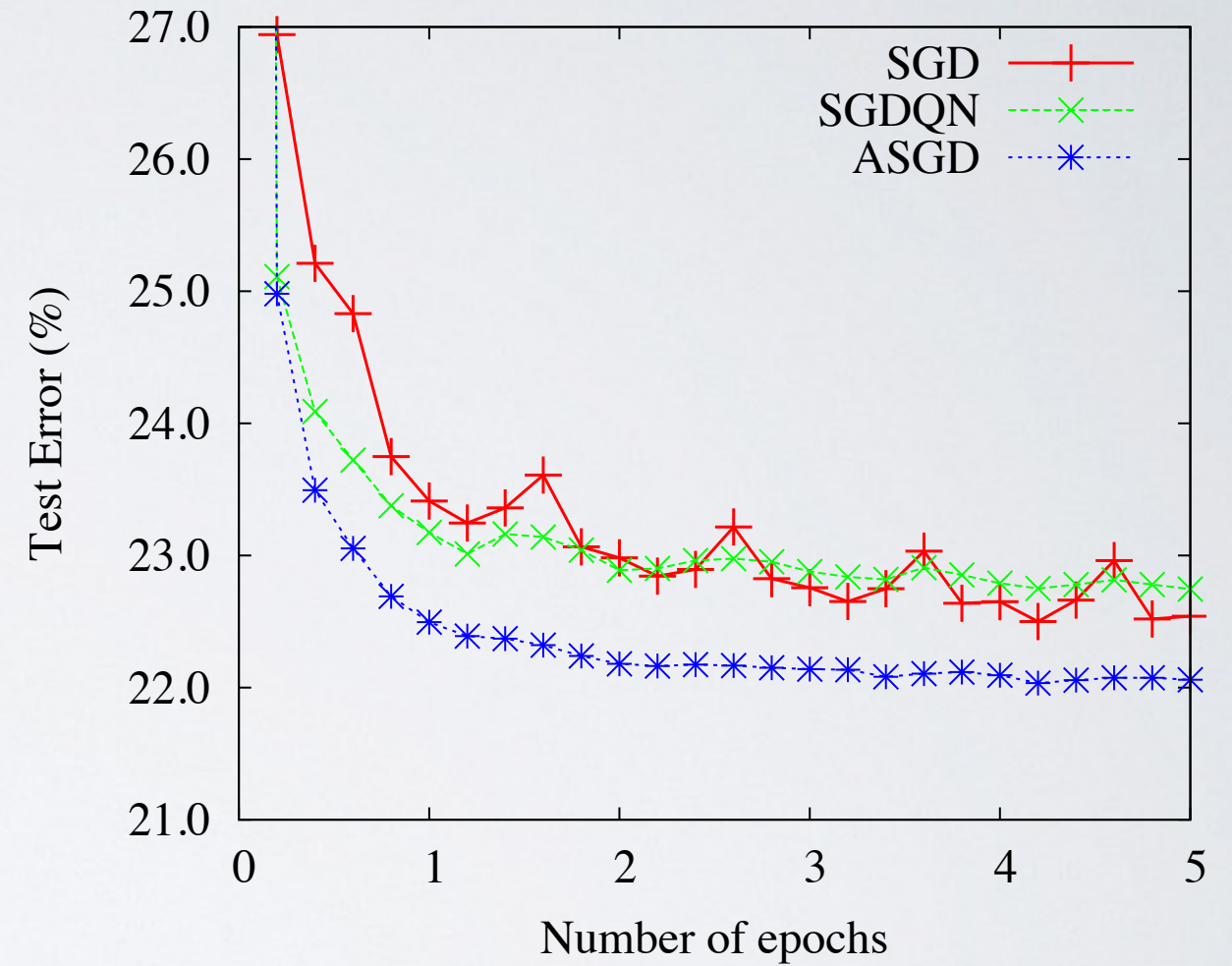
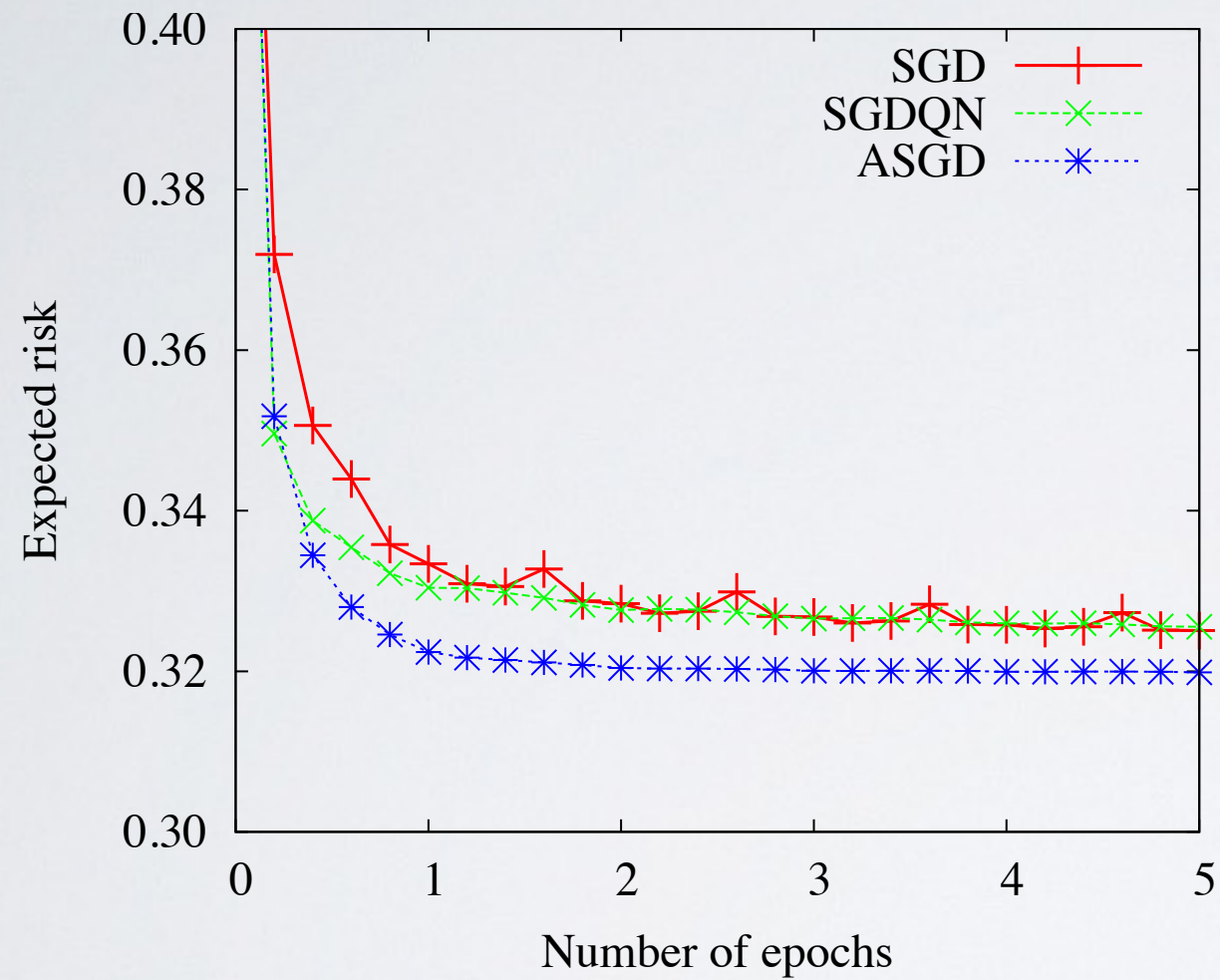
Lasso (Tibshirani, 1996)

$$Q_{\text{lasso}} = \lambda |w|_1 + \frac{1}{2} (y - w^\top \Phi(x))^2$$
$$w = (u_1 - v_1, \dots, u_d - v_d)$$
$$\Phi(x) \in \mathbb{R}^d, \quad y \in \mathbb{R}, \quad \lambda > 0$$

$$u_i \leftarrow [u_i - \gamma_t (\lambda - (y_t - w^\top \Phi(x_t)) \Phi_i(x_t))]_+$$
$$v_i \leftarrow [v_i - \gamma_t (\lambda + (y_t - w_t^\top \Phi(x_t)) \Phi_i(x_t))]_+$$

with notation $[x]_+ = \max\{0, x\}$.

RESULTS



RESULTS

Algorithm	Time	Test Error
<i>Hinge loss SVM, $\lambda = 10^{-4}$.</i>		
SVMLIGHT	23,642 s.	6.02 %
SVMPELF	66 s.	6.03 %
SGD	1.4 s.	6.02 %
<i>Log loss SVM, $\lambda = 10^{-5}$.</i>		
TRON (-e0.01)	30 s.	5.68 %
TRON (-e0.001)	44 s.	5.70 %
SGD	2.3 s.	5.66 %

PRACTICAL CONSIDERATIONS

- Mini - batches
 - reduces variance
- Momentum
 - further reduces variance

SGD FOR DEEP LEARNING

- Faster training
 - a lot of data!

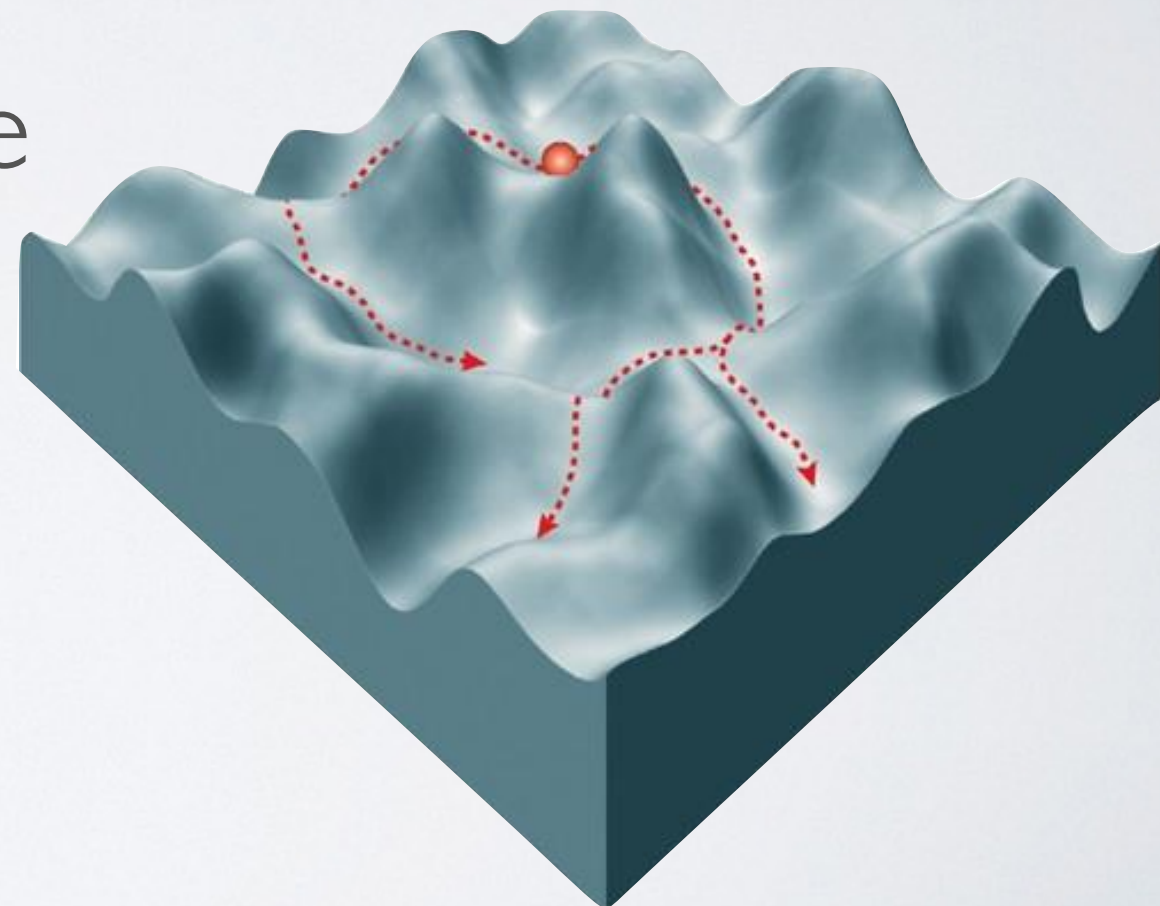


SGD FOR DEEP LEARNING

- Less overfitting
 - shorter training
 - fewer passes over data

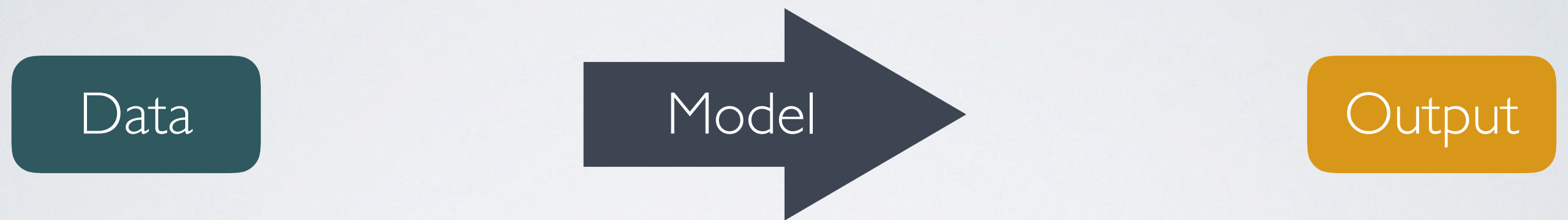
SGD FOR DEEP LEARNING

- Saddle points
 - exponentially more saddle points than local minima
 - GD gets stuck on saddle points



DISCUSSION

PAPER ASSIGNMENT



PAPER ASSIGNMENT

Data

PAPER ASSIGNMENT

Data

S1

P4

P6

P8

...

S2

P3

P6

P7

...

S3

P1

P2

P3

...

S4

P7

P6

P5

...

PAPER ASSIGNMENT

Data

	S1	S2	S3	S4
P1	3	3	1	5
P2	4	2	2	3
P3	1	1	3	2
...

PAPER ASSIGNMENT

Output

PAPER ASSIGNMENT

Output

S1

S2

S3

S4

P4

P3

P1

P7

PAPER ASSIGNMENT

Output

	S1	S2	S3	S4
P1	1	0	0	0
P2	0	0	0	1
P3	0	1	0	0
...				

PAPER ASSIGNMENT

Output

	S1	S2	S3	S4
P1	1	0	0	0
P2	0	0	0	1
P3	0	1	0	0
...				

sum up to 1

PAPER ASSIGNMENT

Output

	S1	S2	S3	S4	
P1	1	0	0	0	sum up to 1
P2	0	0	0	1	
P3	0	1	0	0	
...					

RELAXATION

- Keep summation constraints
- Relax $\{0, 1\}$ to $[0, 1]$

PAPER ASSIGNMENT

Output

	S1	S2	S3	S4	
P1	0.9	0	0	0.1	sum up to 1
P2	0	0.1	0	0.9	
P3	0.1	0.9	0	0	
...					

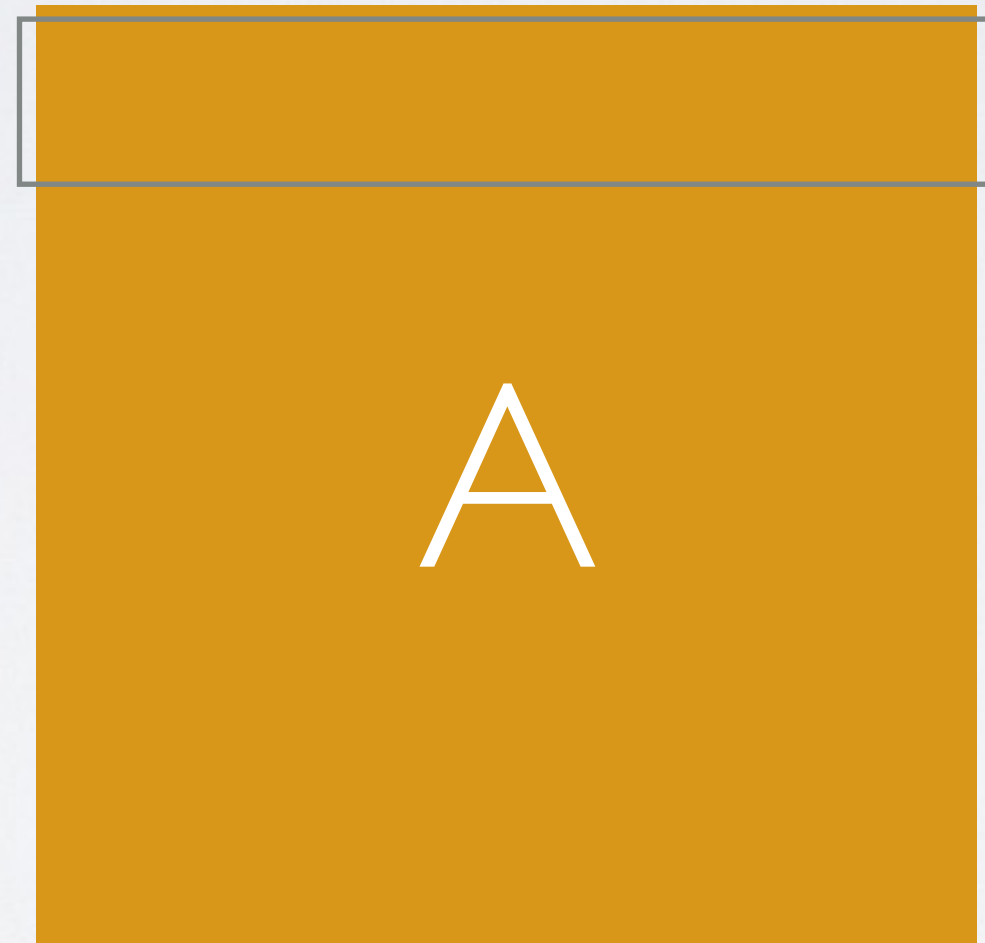
sum up to 1

DOUBLY STOCHASTIC

- Solution 1: Constraints
 - linear program
 - not differentiable
 - how to train?

DOUBLY STOCHASTIC

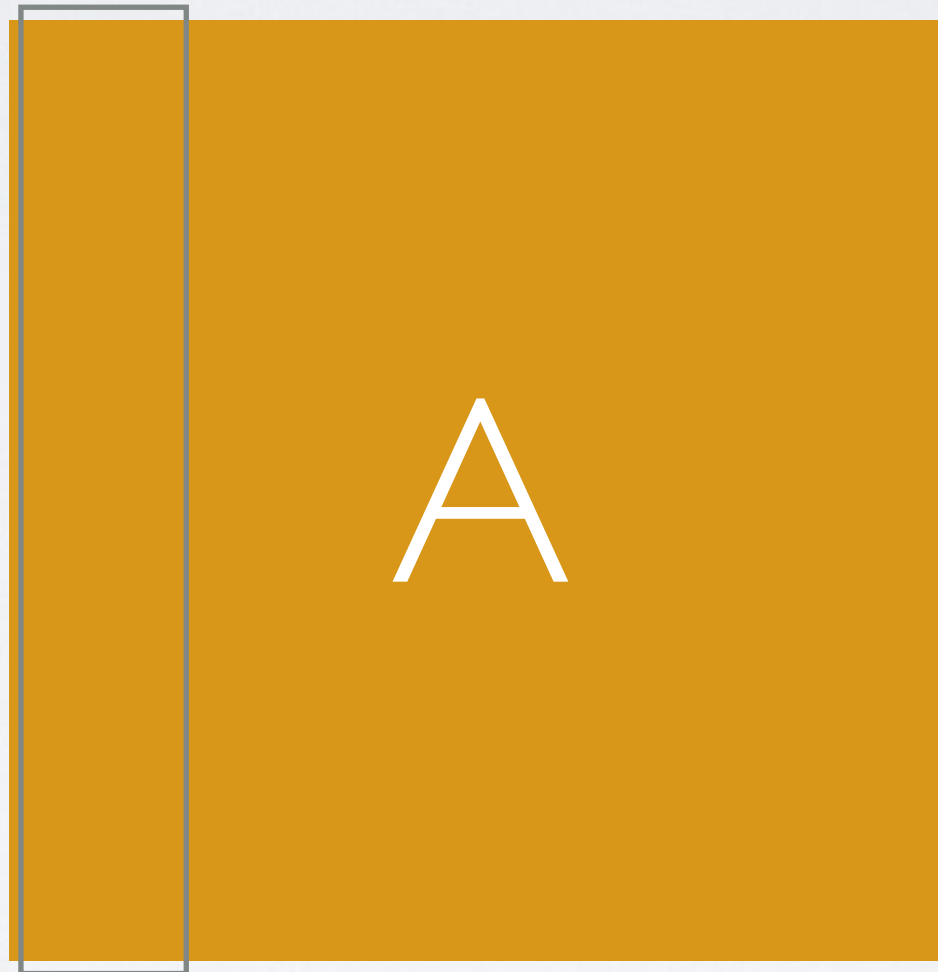
- Solution 2: Approximation



sum and divide

DOUBLY STOCHASTIC

- Solution 2: Approximation



sum and divide

DOUBLY STOCHASTIC

- Solution 2: Approximation
 - Differentiable
 - Trainable with SGD

PAPER ASSIGNMENT

Data

	S1	S2	S3	S4
P1	3	3	1	5
P2	4	2	2	3
P3	1	1	3	2
...

Input matrix D

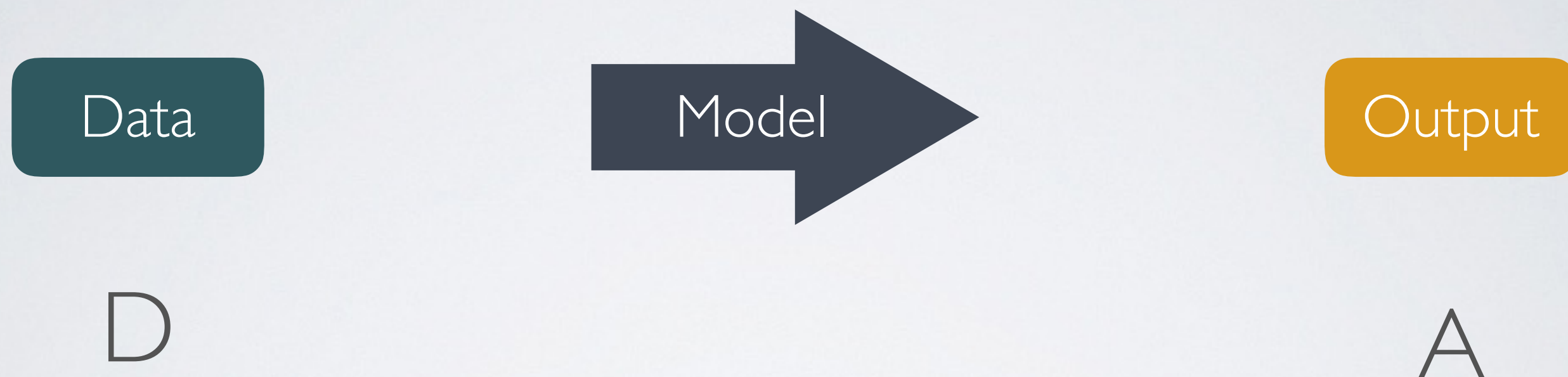
PAPER ASSIGNMENT

Output

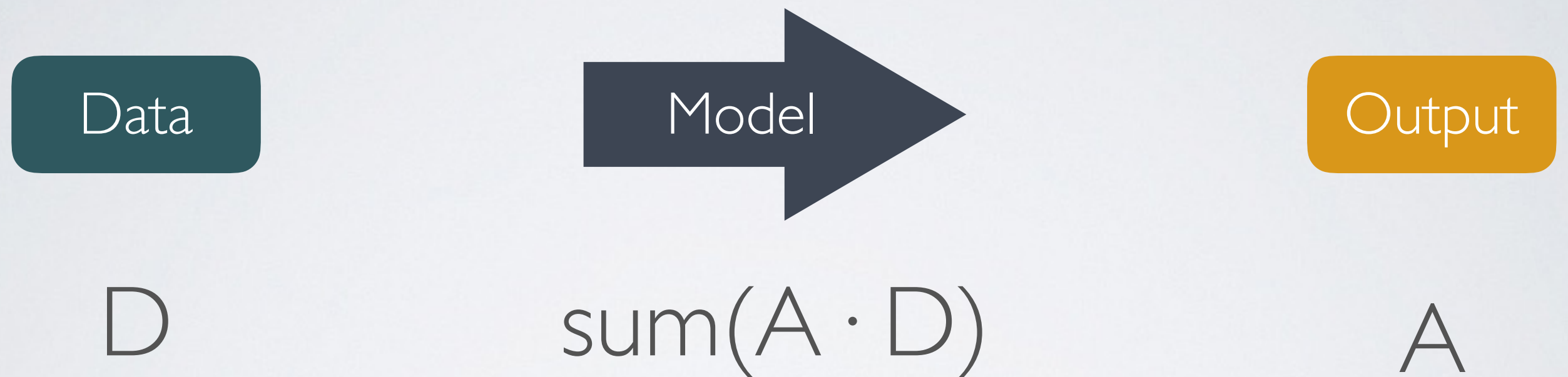
	S1	S2	S3	S4
P1	0.9	0	0	0.1
P2	0	0.1	0	0.9
P3	0.1	0.9	0	0
...				

Output matrix A

TRAINING LOSS



TRAINING LOSS



TENSORFLOW IMPLEMENTATION