

TTIC 31190: Natural Language Processing

Kevin Gimpel
Spring 2018

Lecture 6: Learning for Classification; Language Modeling

- assignment 1 due today
- questions?
- if you want to use late day(s), state that on your report

- assignment 2 will be posted tomorrow
- start thinking about your project, who you might want to work with, etc.

- short quiz at start of class Wed., April 18th
- covering material up to and including Mon., April 9th
- don't stress about it
- grading will be check-minus/check/check-plus

Roadmap

- words, morphology, lexical semantics
- text classification
- simple neural methods for NLP
- language modeling and word embeddings
- recurrent/recursive/convolutional networks in NLP
- sequence labeling, HMMs, dynamic programming
- syntax and syntactic parsing
- semantics, compositionality, semantic parsing
- machine translation and other NLP tasks

Text Classification

- datasets
- classification
 - modeling
 - inference
 - learning

Classifiers

- one simple type:
 - for any input \mathbf{x} , assign a score to each label y

$$\text{score}(\mathbf{x}, y, \mathbf{w})$$

- classify by choosing highest-scoring label:

$$\text{classify}(\mathbf{x}, \mathbf{w}) = \underset{y}{\operatorname{argmax}} \text{score}(\mathbf{x}, y, \mathbf{w})$$

Linear Models

- parameters are arranged in a vector \mathbf{w}
- score function is linear in \mathbf{w} :

$$\text{score}(\mathbf{x}, y, \mathbf{w}) = \mathbf{w}^\top \mathbf{f}(\mathbf{x}, y) = \sum_i w_i f_i(\mathbf{x}, y)$$

- \mathbf{f} : vector of feature functions

Modeling, Inference, Learning

inference: solve argmax

modeling: define score function

$$\operatorname{classify}(\boldsymbol{x}, \boldsymbol{w}) = \operatorname{argmax}_y \operatorname{score}(\boldsymbol{x}, y, \boldsymbol{w})$$

learning: choose \boldsymbol{w}

- **Learning:** How do we choose the weights \boldsymbol{w} ?

Regularized Empirical Risk Minimization

$$\hat{\boldsymbol{w}} = \operatorname{argmin}_{\boldsymbol{w}} \sum_{i=1}^{|\mathcal{T}|} \operatorname{loss}(\boldsymbol{x}^{(i)}, y^{(i)}, \boldsymbol{w}) + \underbrace{\lambda R(\boldsymbol{w})}_{\text{regularization term}}$$

Diagram illustrating the components of the regularized empirical risk minimization equation:


- surrogate loss**: Points to the sum of losses $\sum_{i=1}^{|\mathcal{T}|} \operatorname{loss}(\boldsymbol{x}^{(i)}, y^{(i)}, \boldsymbol{w})$.
- regularization strength**: Points to the regularization parameter λ .
- regularization term**: Points to the regularization function $R(\boldsymbol{w})$.

Regularized Empirical Risk Minimization

$$\hat{\boldsymbol{w}} = \operatorname{argmin}_{\boldsymbol{w}} \sum_{i=1}^{|\mathcal{T}|} \operatorname{loss}(\boldsymbol{x}^{(i)}, y^{(i)}, \boldsymbol{w}) + \lambda R(\boldsymbol{w})$$



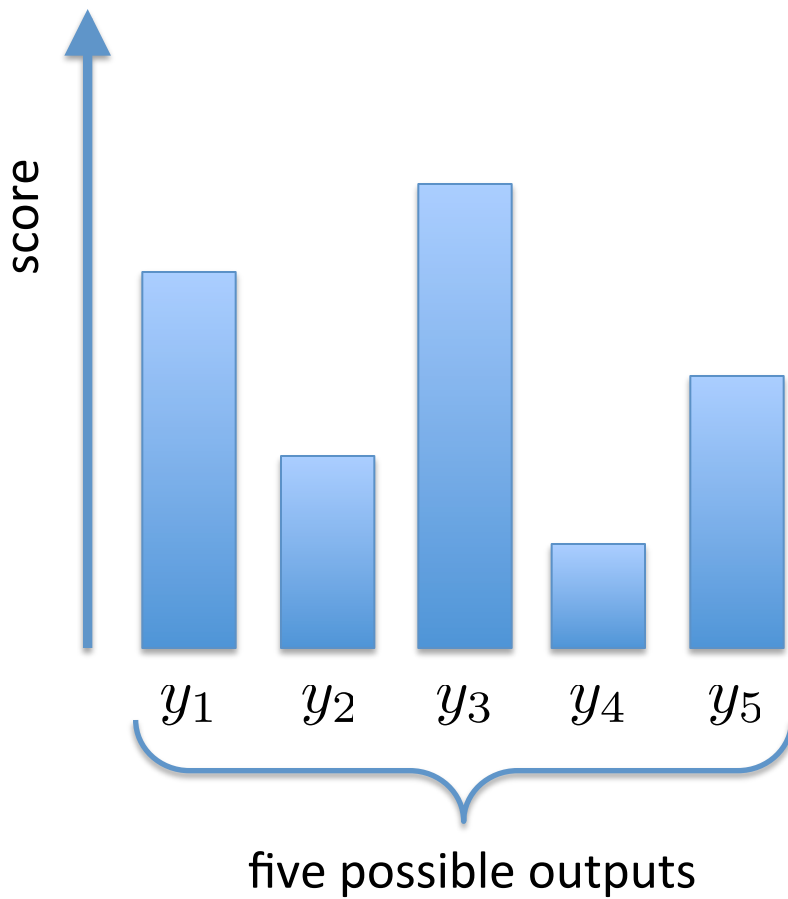
**encourages model
to fit the training
data well**



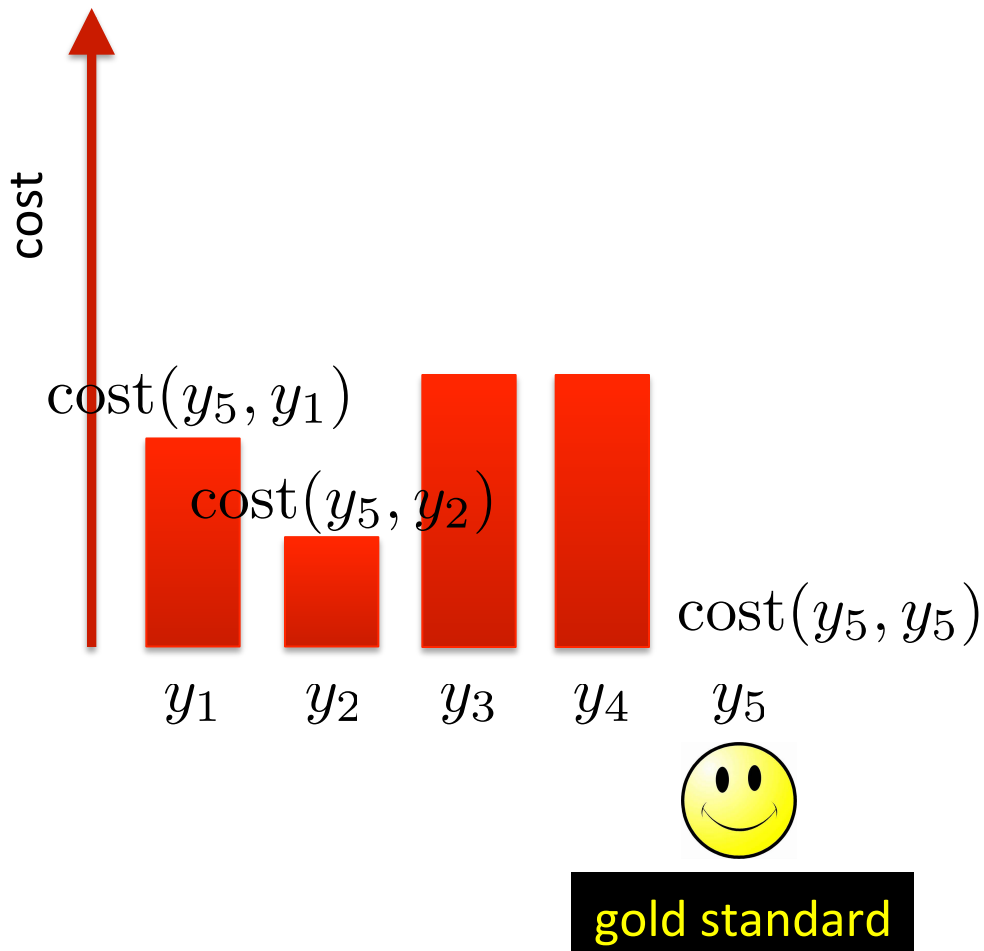
**encourages model to be
“simpler” in the hope that
this will help it to
generalize to new data**

Visualization

for a single input x

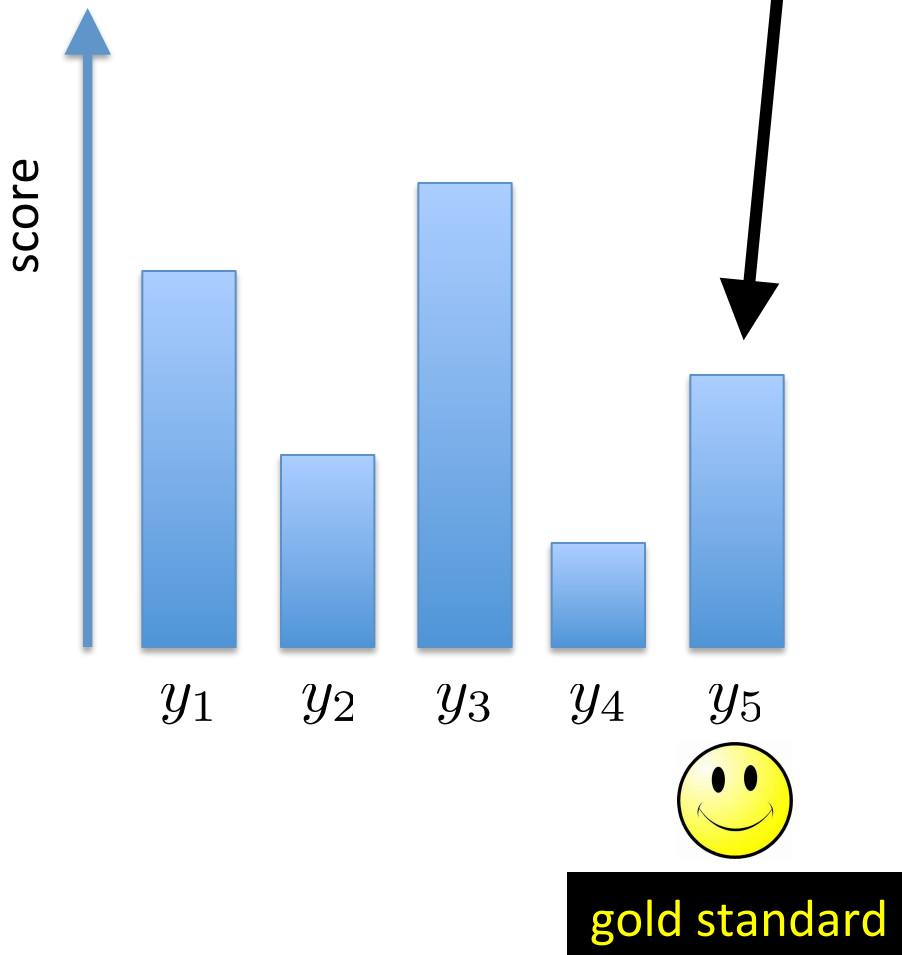


Visualization for a single input x



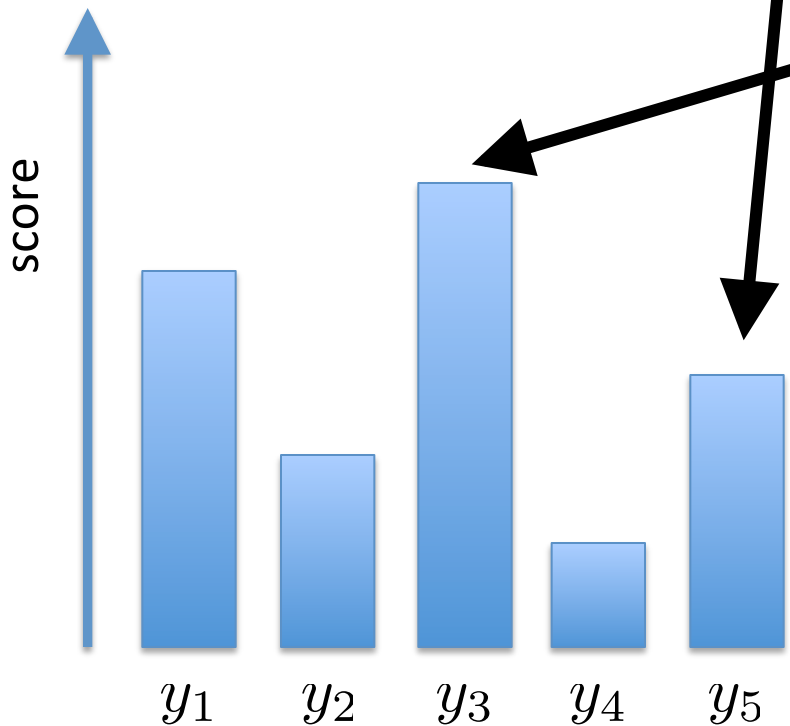
perceptron loss:

$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w}) + \underbrace{\max_{y' \in \mathcal{L}} \text{score}(\mathbf{x}, y', \mathbf{w})}$$



perceptron loss:

$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w}) + \max_{y' \in \mathcal{L}} \text{score}(\mathbf{x}, y', \mathbf{w})$$

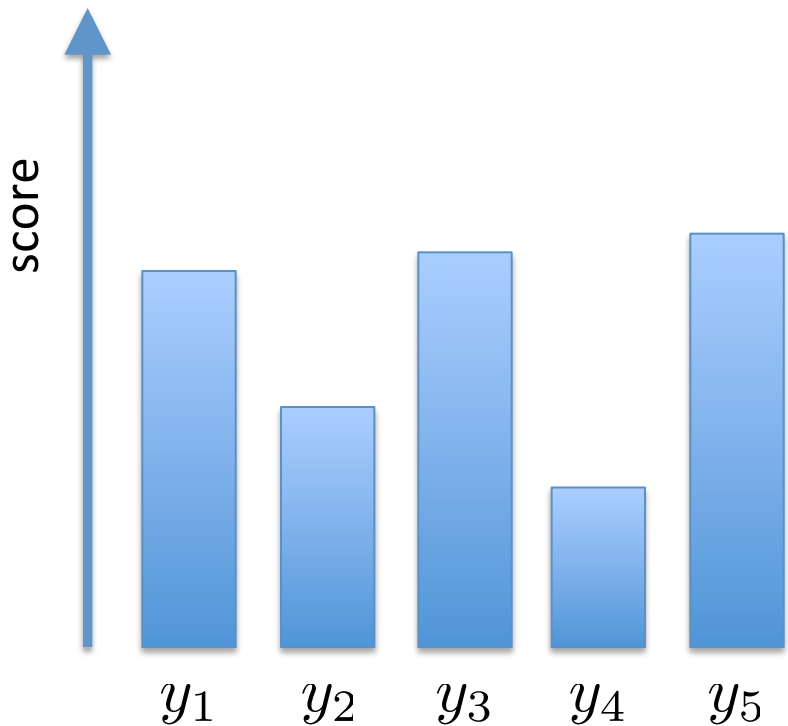


gold standard

effect of learning?

perceptron loss:

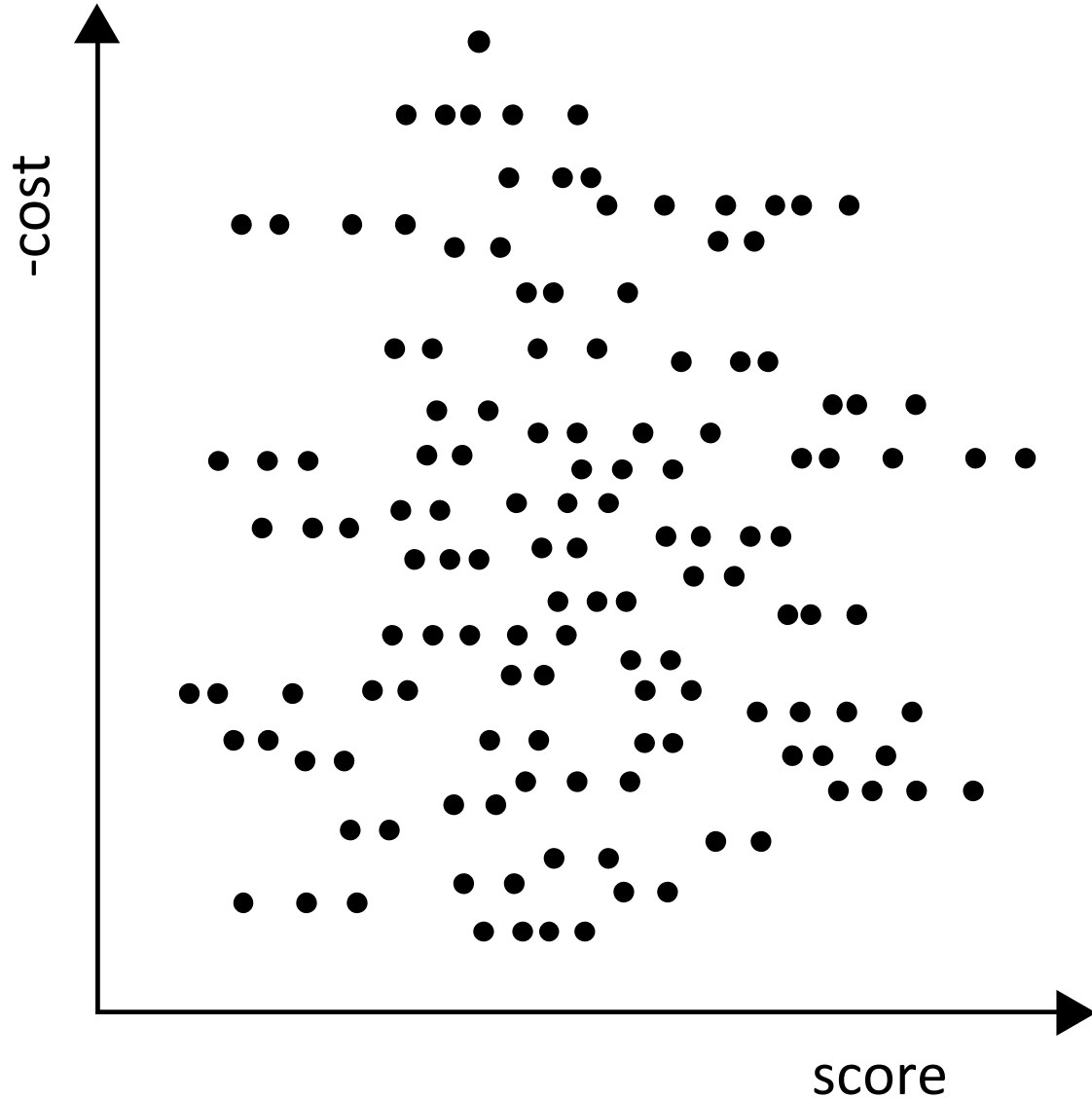
$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w}) + \max_{y' \in \mathcal{L}} \text{score}(\mathbf{x}, y', \mathbf{w})$$



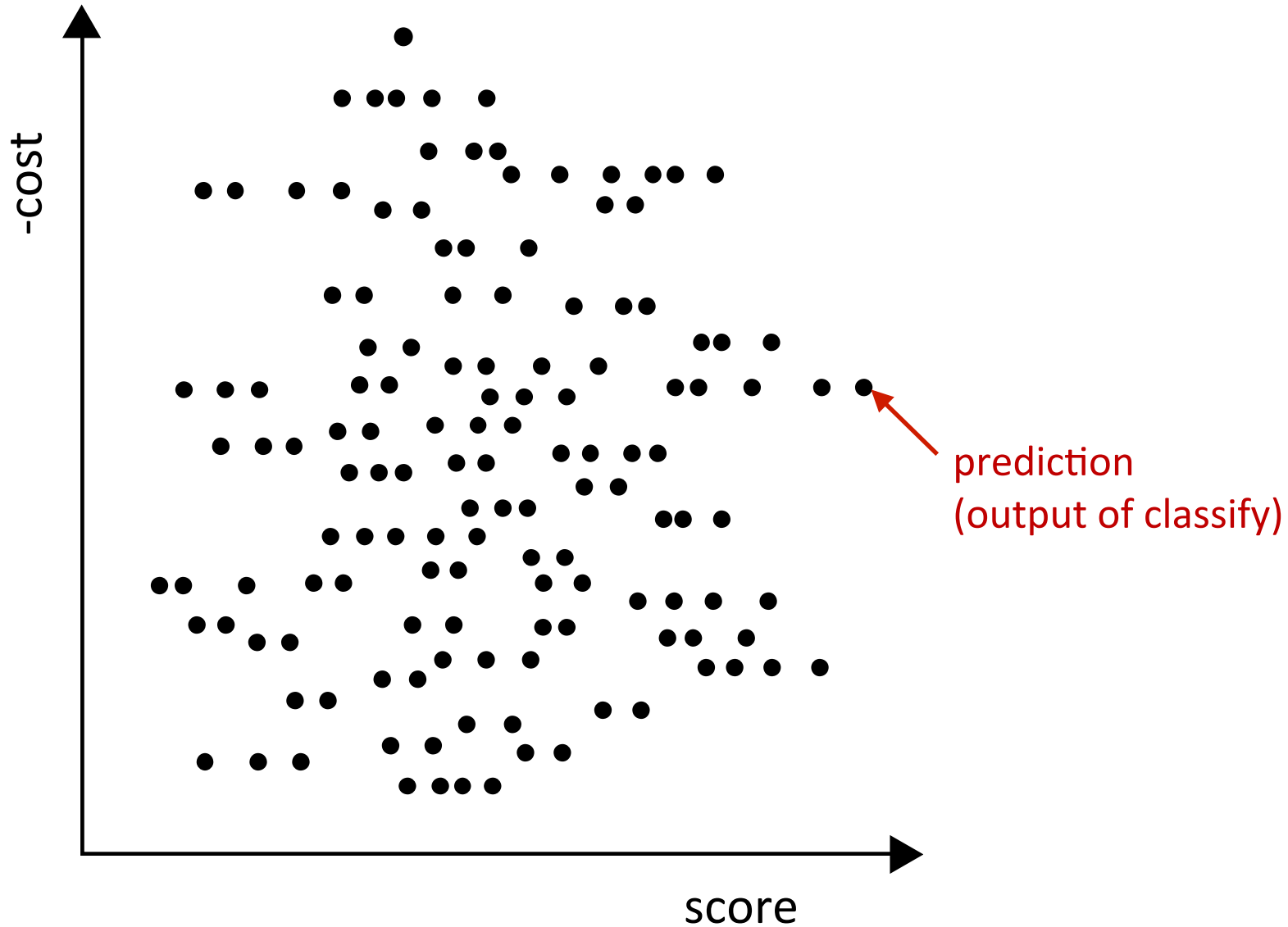
gold standard

effect of learning:
gold standard will
have highest score

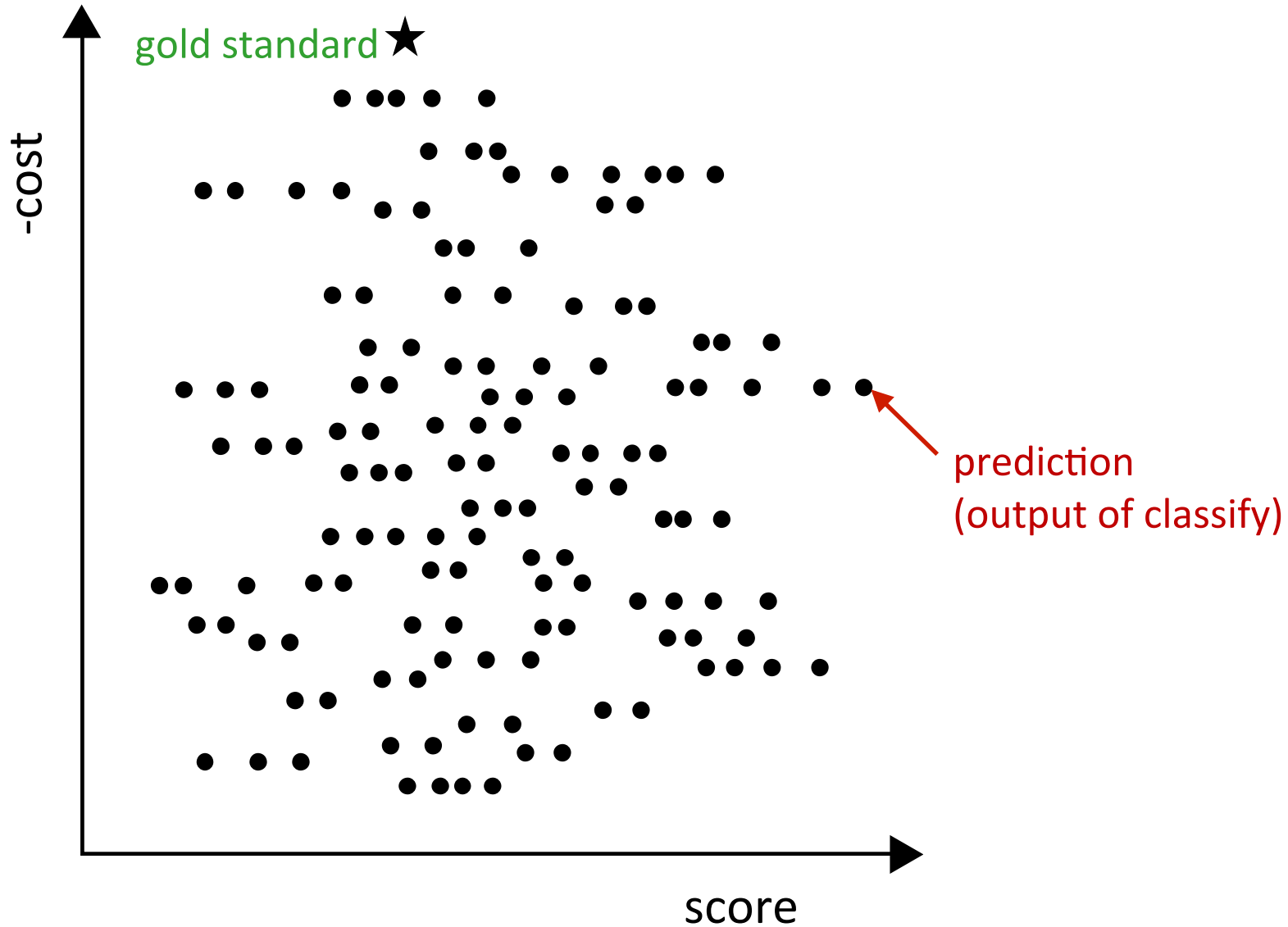
Visualization for a single x :
each point is a possible y

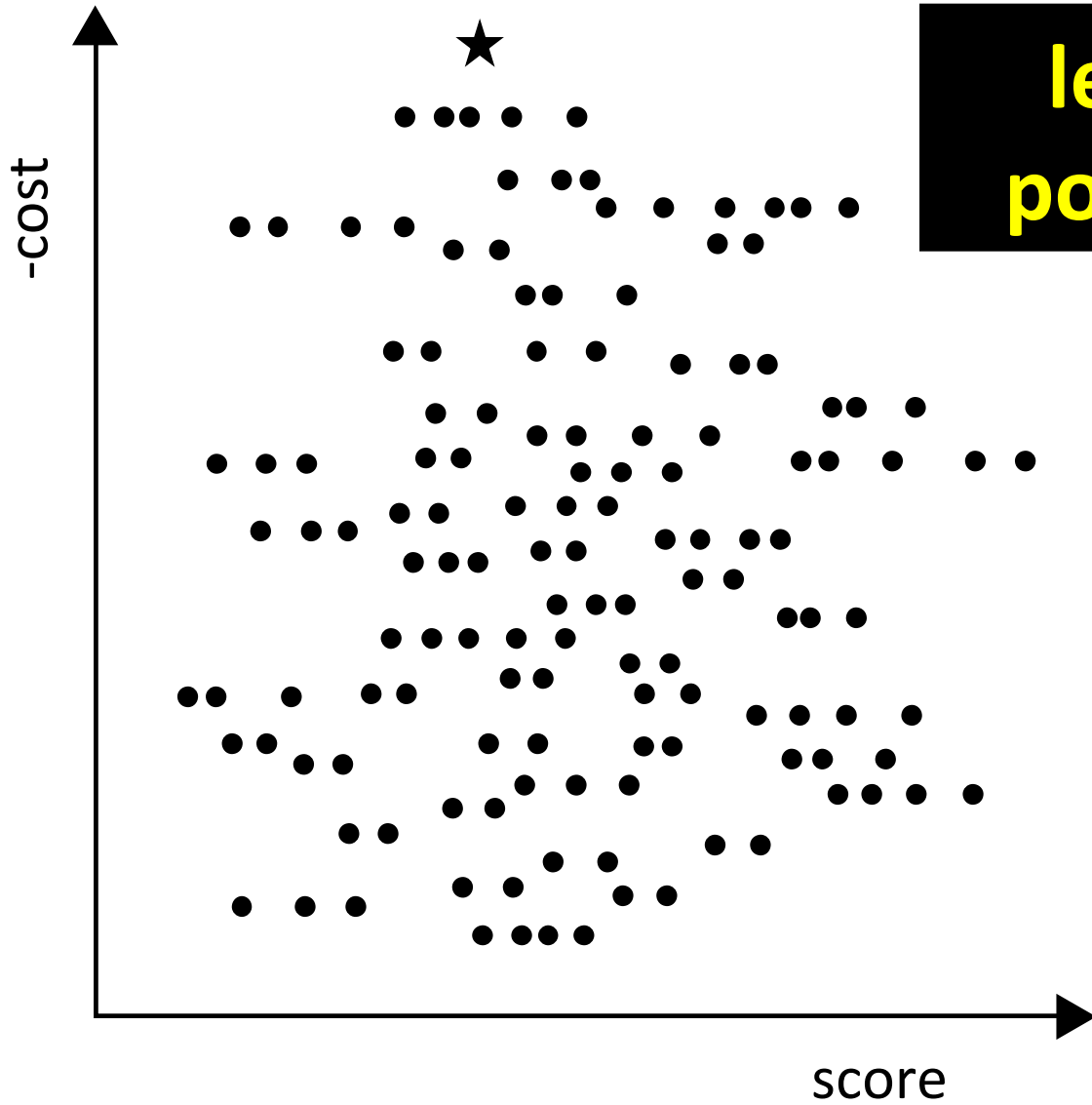


Visualization for a single x :
each point is a possible y

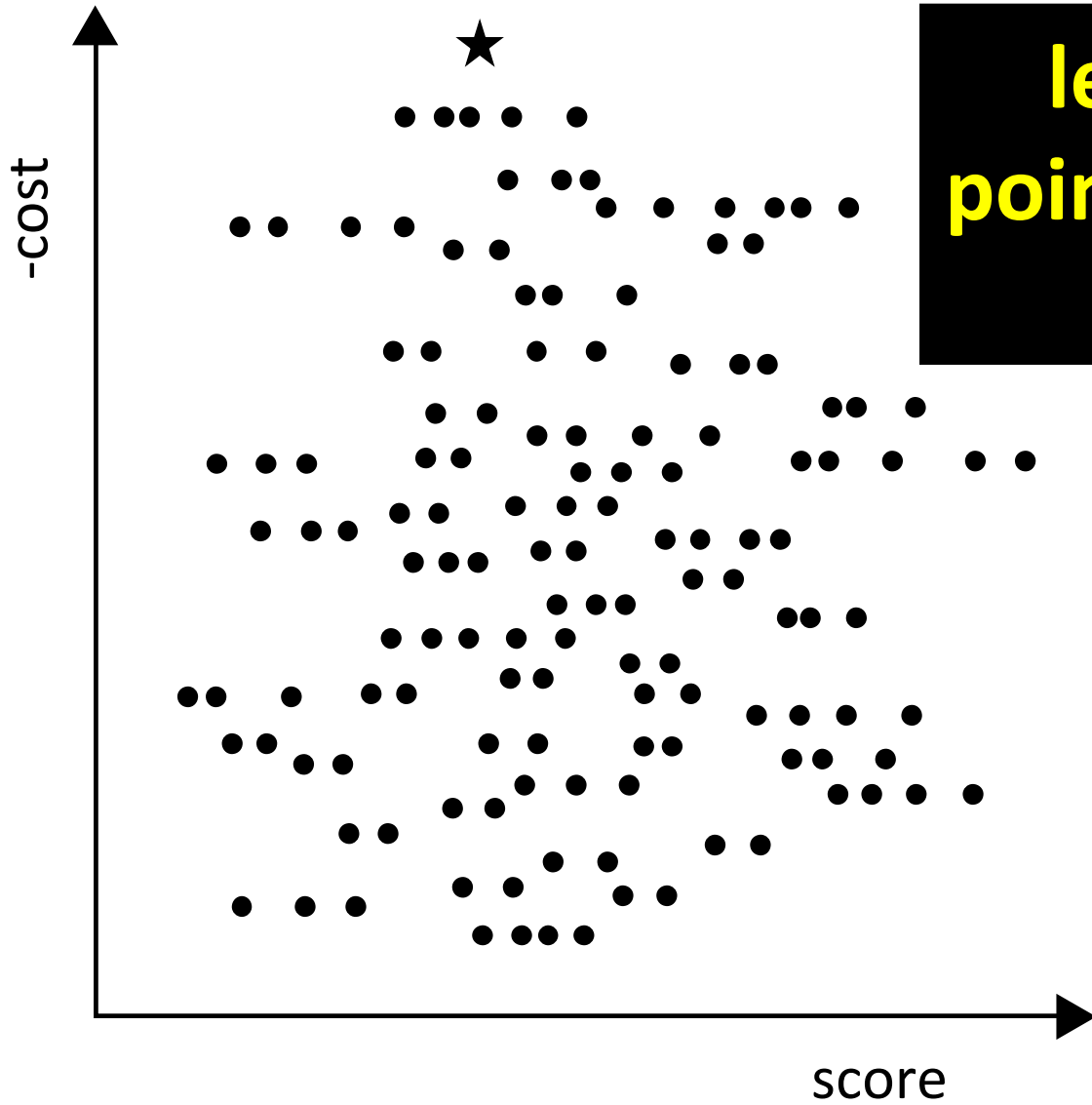


Visualization for a single x : each point is a possible y

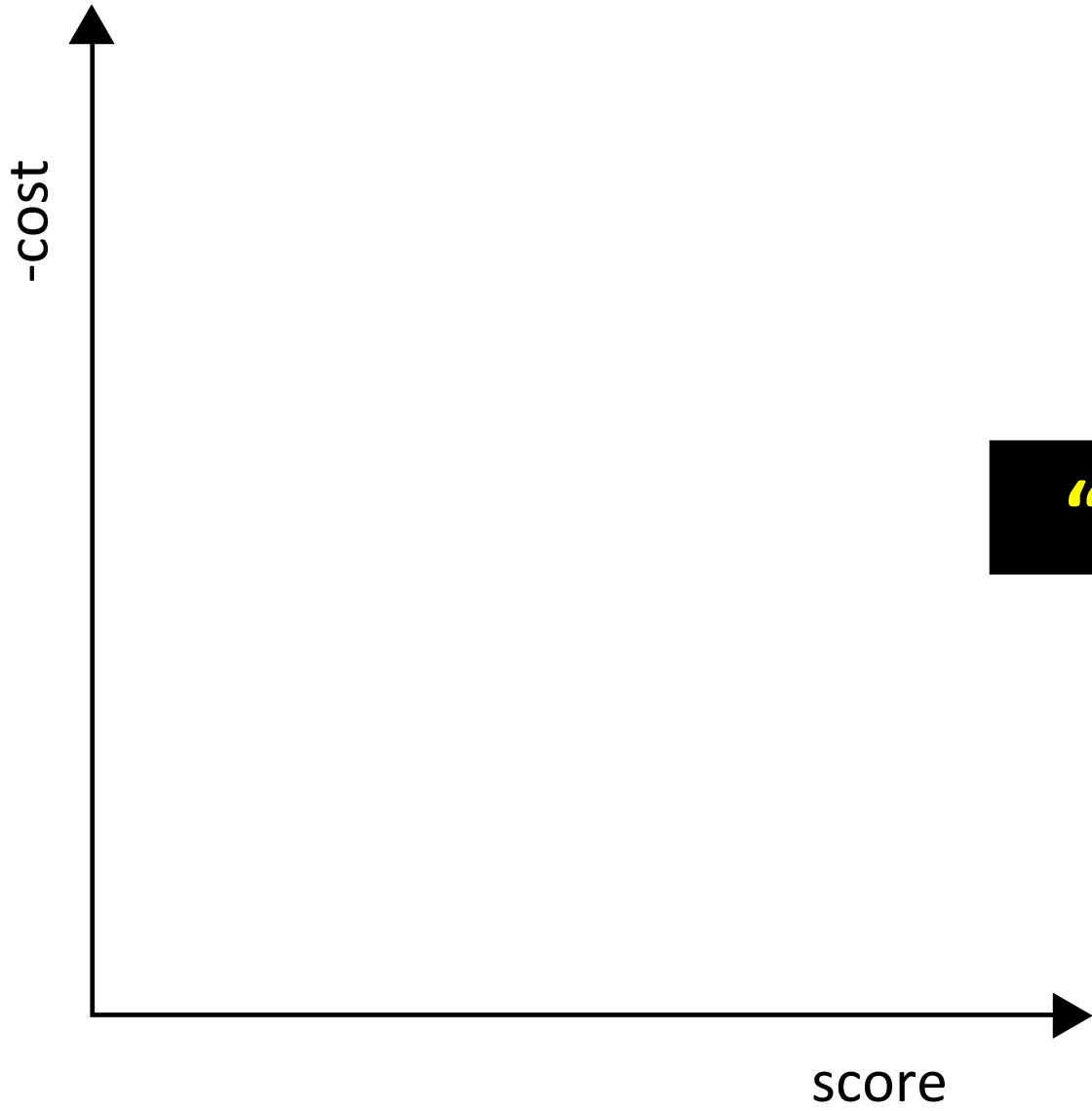




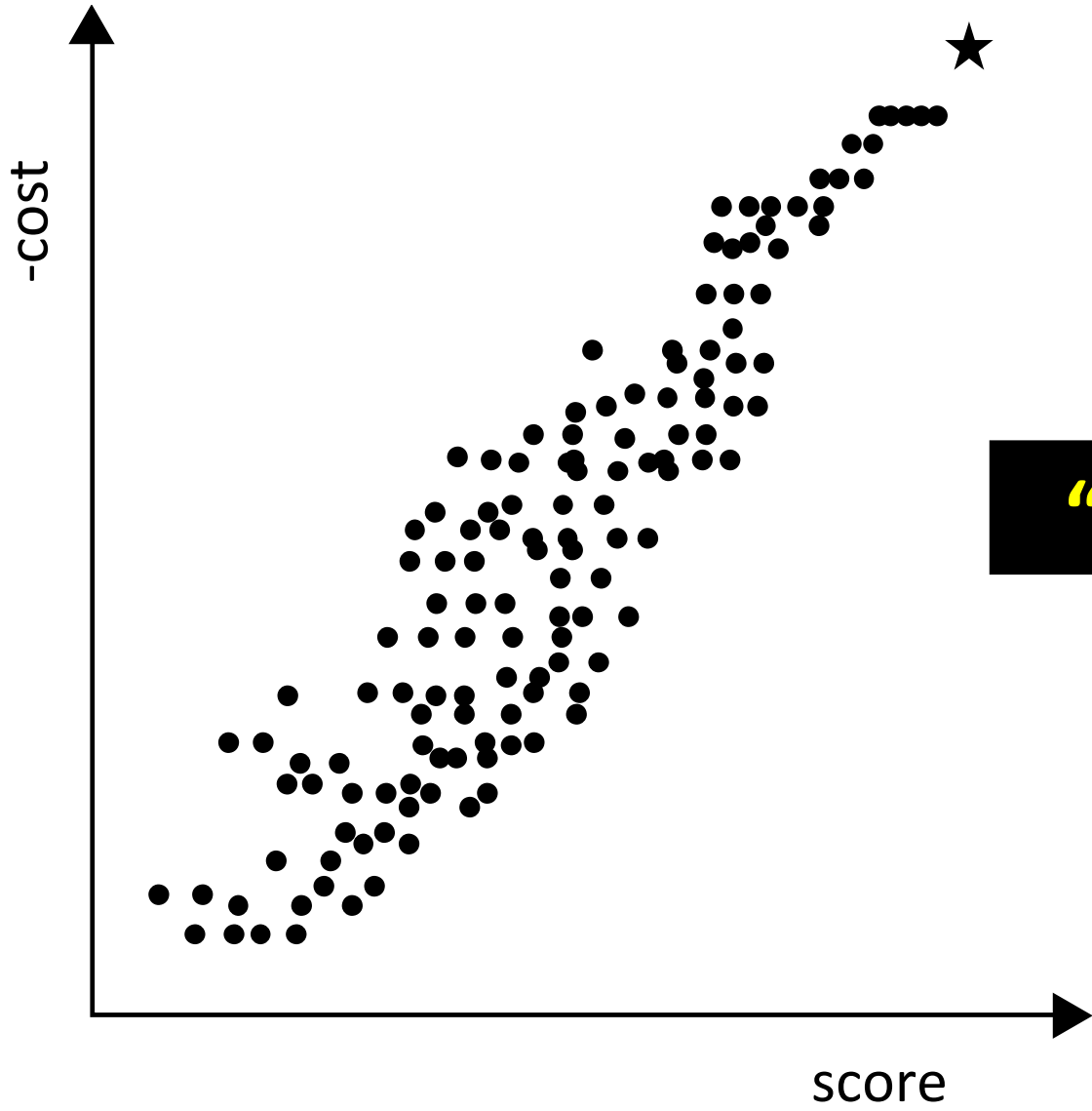
**learning moves
points in this plot**



**learning moves
points left or right in
this plot**

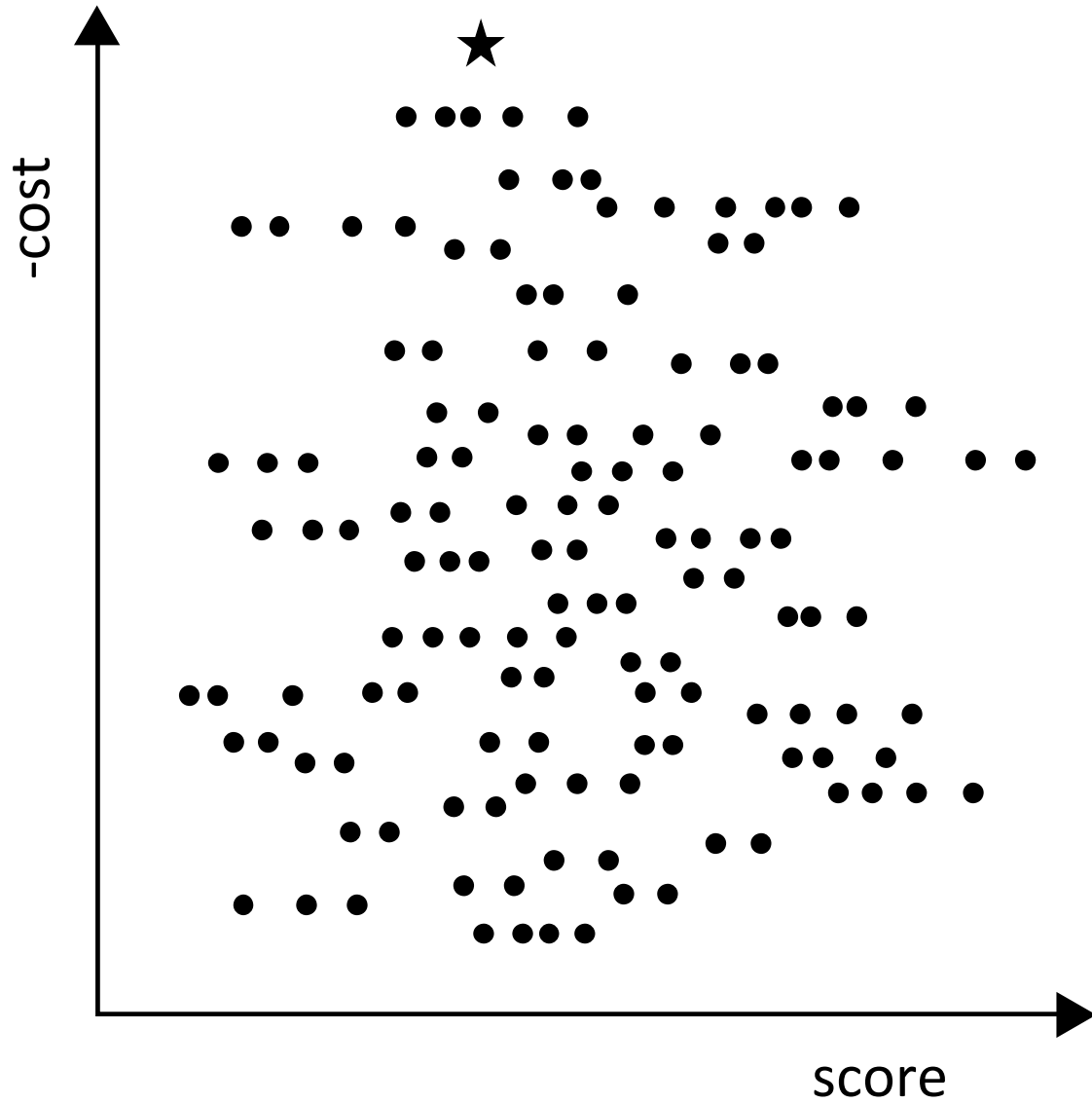


“ideal” model?

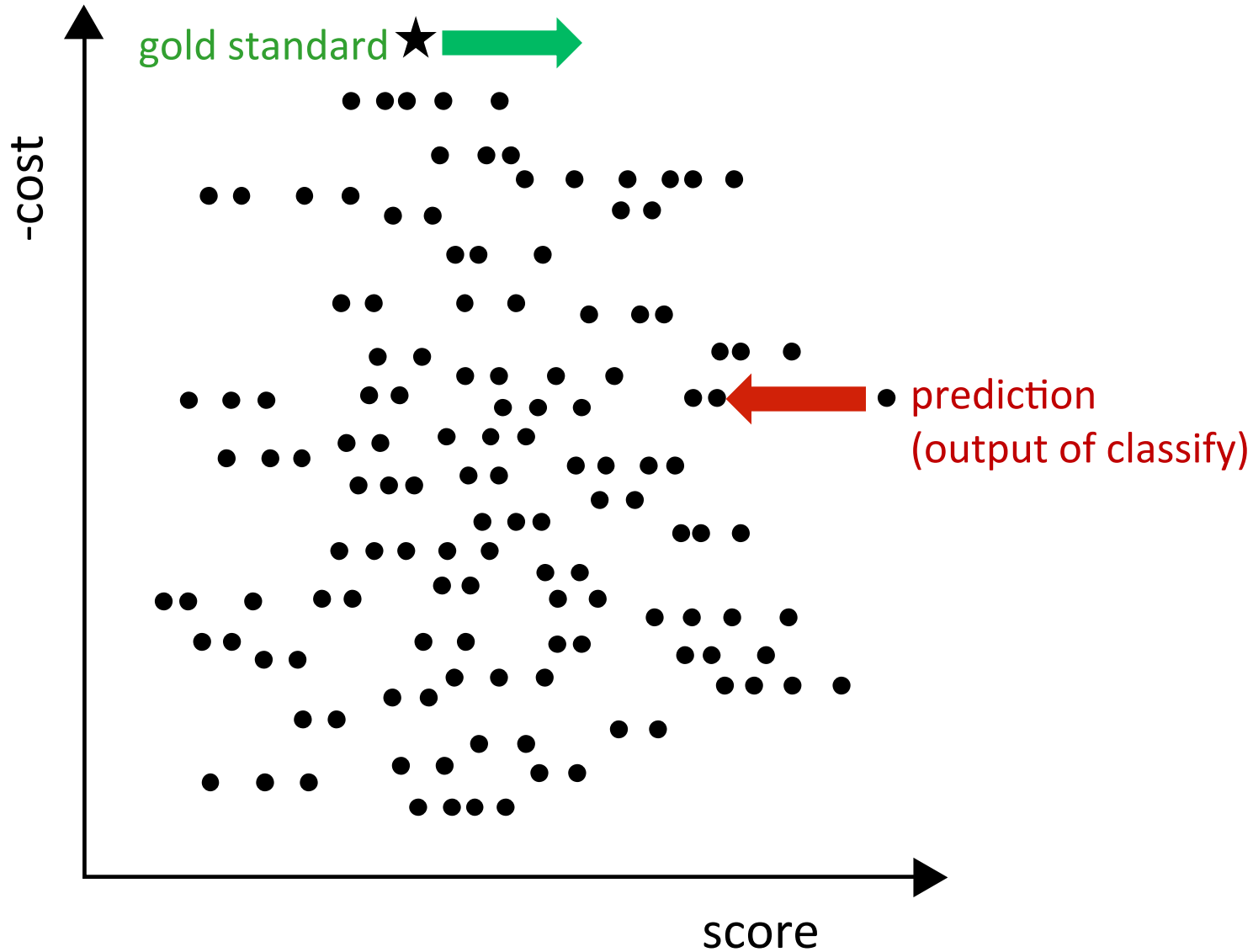


“ideal” model?

Perceptron Loss?



Perceptron Loss



Losses for Linear Models

$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w}) + \max_{y' \in \mathcal{L}} \text{score}(\mathbf{x}, y', \mathbf{w})$$



$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = -\mathbf{w}^\top \mathbf{f}(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, y')$$

Losses for Linear Models

$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w}) + \max_{y' \in \mathcal{L}} \text{score}(\mathbf{x}, y', \mathbf{w})$$



$$\begin{aligned} \text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) &= -\mathbf{w}^\top \mathbf{f}(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, y') \\ &= -\sum_i w_i f_i(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \sum_i w_i f_i(\mathbf{x}, y') \end{aligned}$$

Loss Subgradients for Linear Models

- some of our loss functions are not differentiable:

$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = - \sum_i w_i f_i(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \sum_i w_i f_i(\mathbf{x}, y')$$

- but they are subdifferentiable:

$$\frac{\partial \text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w})}{\partial w_j} =$$

Loss Subgradients for Linear Models

- some of our loss functions are not differentiable:

$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = - \sum_i w_i f_i(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \sum_i w_i f_i(\mathbf{x}, y')$$

- but they are subdifferentiable:

$$\frac{\partial \text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w})}{\partial w_j} = -f_j(\mathbf{x}, y) +$$

Loss Subgradients for Linear Models

- some of our loss functions are not differentiable:

$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = - \sum_i w_i f_i(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \sum_i w_i f_i(\mathbf{x}, y')$$

- but they are subdifferentiable:

$$\frac{\partial \text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w})}{\partial w_j} = -f_j(\mathbf{x}, y) + \frac{\partial}{\partial w_j} \max_{y' \in \mathcal{L}} \sum_i w_i f_i(\mathbf{x}, y')$$

Loss Subgradients for Linear Models

- some of our loss functions are not differentiable:

$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = - \sum_i w_i f_i(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \sum_i w_i f_i(\mathbf{x}, y')$$

- but they are subdifferentiable:

$$\frac{\partial}{\partial w_j} \max_{y' \in \mathcal{L}} \sum_i w_i f_i(\mathbf{x}, y') =$$



find subgradient of the function
that achieves the max

Loss Subgradients for Linear Models

- some of our loss functions are not differentiable:

$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = - \sum_i w_i f_i(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \sum_i w_i f_i(\mathbf{x}, y')$$

- but they are subdifferentiable:

$$\frac{\partial}{\partial w_j} \max_{y' \in \mathcal{L}} \sum_i w_i f_i(\mathbf{x}, y') = \frac{\partial}{\partial w_j} \sum_i w_i f_i(\mathbf{x}, \text{classify}(\mathbf{x}, \mathbf{w}))$$



find subgradient of the function
that achieves the max

Loss Subgradients for Linear Models

- some of our loss functions are not differentiable:

$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = - \sum_i w_i f_i(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \sum_i w_i f_i(\mathbf{x}, y')$$

- but they are subdifferentiable:

$$\frac{\partial}{\partial w_j} \max_{y' \in \mathcal{L}} \sum_i w_i f_i(\mathbf{x}, y') = \frac{\partial}{\partial w_j} \sum_i w_i f_i(\mathbf{x}, \text{classify}(\mathbf{x}, \mathbf{w})) = f_j(\mathbf{x}, \text{classify}(\mathbf{x}, \mathbf{w}))$$

find subgradient of the function
that achieves the max

Loss Subgradients for Linear Models

- perceptron loss:

$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = - \sum_i w_i f_i(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \sum_i w_i f_i(\mathbf{x}, y')$$

- subderivative for a single parameter:

$$\frac{\partial \text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w})}{\partial w_j} = -f_j(\mathbf{x}, y) + f_j(\mathbf{x}, \text{classify}(\mathbf{x}, \mathbf{w}))$$

Loss Subgradients for Linear Models

- perceptron loss and subgradient:

$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = - \sum_i w_i f_i(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \sum_i w_i f_i(\mathbf{x}, y')$$

$$\frac{\partial \text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w})}{\partial w_j} = -f_j(\mathbf{x}, y) + f_j(\mathbf{x}, \text{classify}(\mathbf{x}, \mathbf{w}))$$

- parameter update:

$$w_j \leftarrow w_j + f_j(\mathbf{x}, y) - f_j(\mathbf{x}, \text{classify}(\mathbf{x}, \mathbf{w}))$$

Loss Subgradients for Linear Models

- perceptron loss and subgradient:

$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = - \sum_i w_i f_i(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \sum_i w_i f_i(\mathbf{x}, y')$$

$$\frac{\partial \text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w})}{\partial w_j} = -f_j(\mathbf{x}, y) + f_j(\mathbf{x}, \text{classify}(\mathbf{x}, \mathbf{w}))$$

- parameter update:

$$w_j \leftarrow w_j + \eta (f_j(\mathbf{x}, y) - f_j(\mathbf{x}, \text{classify}(\mathbf{x}, \mathbf{w})))$$



step size / learning rate for
stochastic subgradient descent

Perceptron Loss with Regularization

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^{|\mathcal{T}|} \operatorname{loss}_{\text{perc}}(\mathbf{x}^{(i)}, y^{(i)}, \mathbf{w}) + \lambda R_{L2}(\mathbf{w})$$

$$R_{L2}(\mathbf{w}) = \|\mathbf{w}\|_2^2 = \sum_i w_i^2$$

- update rule from before:

$$w_j \leftarrow w_j + \eta (f_j(\mathbf{x}, y) - f_j(\mathbf{x}, \operatorname{classify}(\mathbf{x}, \mathbf{w})))$$

- with L2 regularization:

Perceptron Loss with Regularization

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^{|\mathcal{T}|} \operatorname{loss}_{\text{perc}}(\mathbf{x}^{(i)}, y^{(i)}, \mathbf{w}) + \lambda R_{L2}(\mathbf{w})$$

$$R_{L2}(\mathbf{w}) = \|\mathbf{w}\|_2^2 = \sum_i w_i^2$$

- update rule from before:

$$w_j \leftarrow w_j + \eta (f_j(\mathbf{x}, y) - f_j(\mathbf{x}, \operatorname{classify}(\mathbf{x}, \mathbf{w})))$$

- with L2 regularization:

$$w_j \leftarrow w_j + \eta (f_j(\mathbf{x}, y) - f_j(\mathbf{x}, \operatorname{classify}(\mathbf{x}, \mathbf{w})) - 2\lambda w_j)$$

Perceptron Loss with Regularization

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^{|\mathcal{T}|} \operatorname{loss}_{\text{perc}}(\mathbf{x}^{(i)}, y^{(i)}, \mathbf{w}) + \lambda R_{L2}(\mathbf{w})$$

$$R_{L2}(\mathbf{w}) = \|\mathbf{w}\|_2^2 = \sum_i w_i^2$$

- update rule from before:

$$w_j \leftarrow w_j + \eta (f_j(\mathbf{x}, y) - f_j(\mathbf{x}, \operatorname{classify}(\mathbf{x}, \mathbf{w})))$$

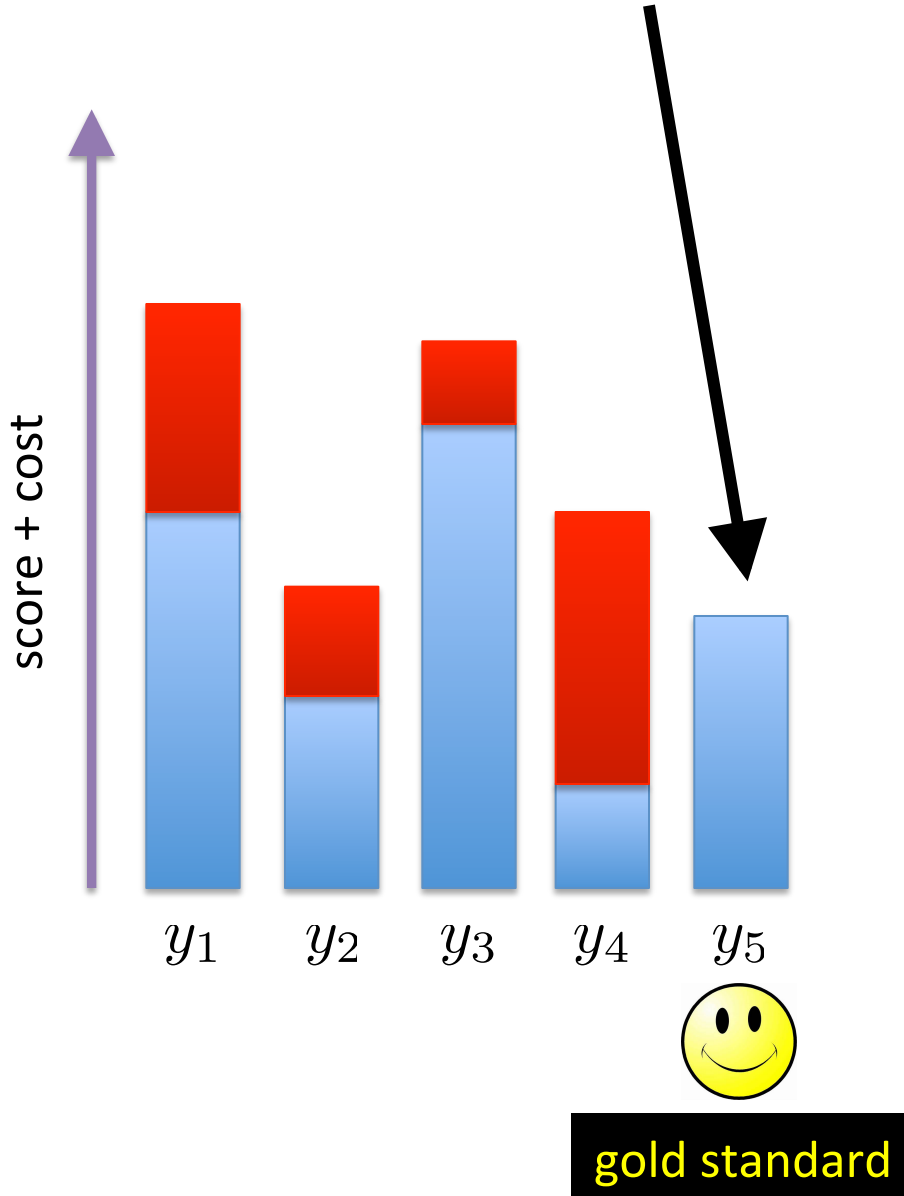
- with L2 regularization:

$$w_j \leftarrow w_j + \eta (f_j(\mathbf{x}, y) - f_j(\mathbf{x}, \operatorname{classify}(\mathbf{x}, \mathbf{w})) - 2\lambda w_j)$$

pushes weights closer to zero (“weight decay”)

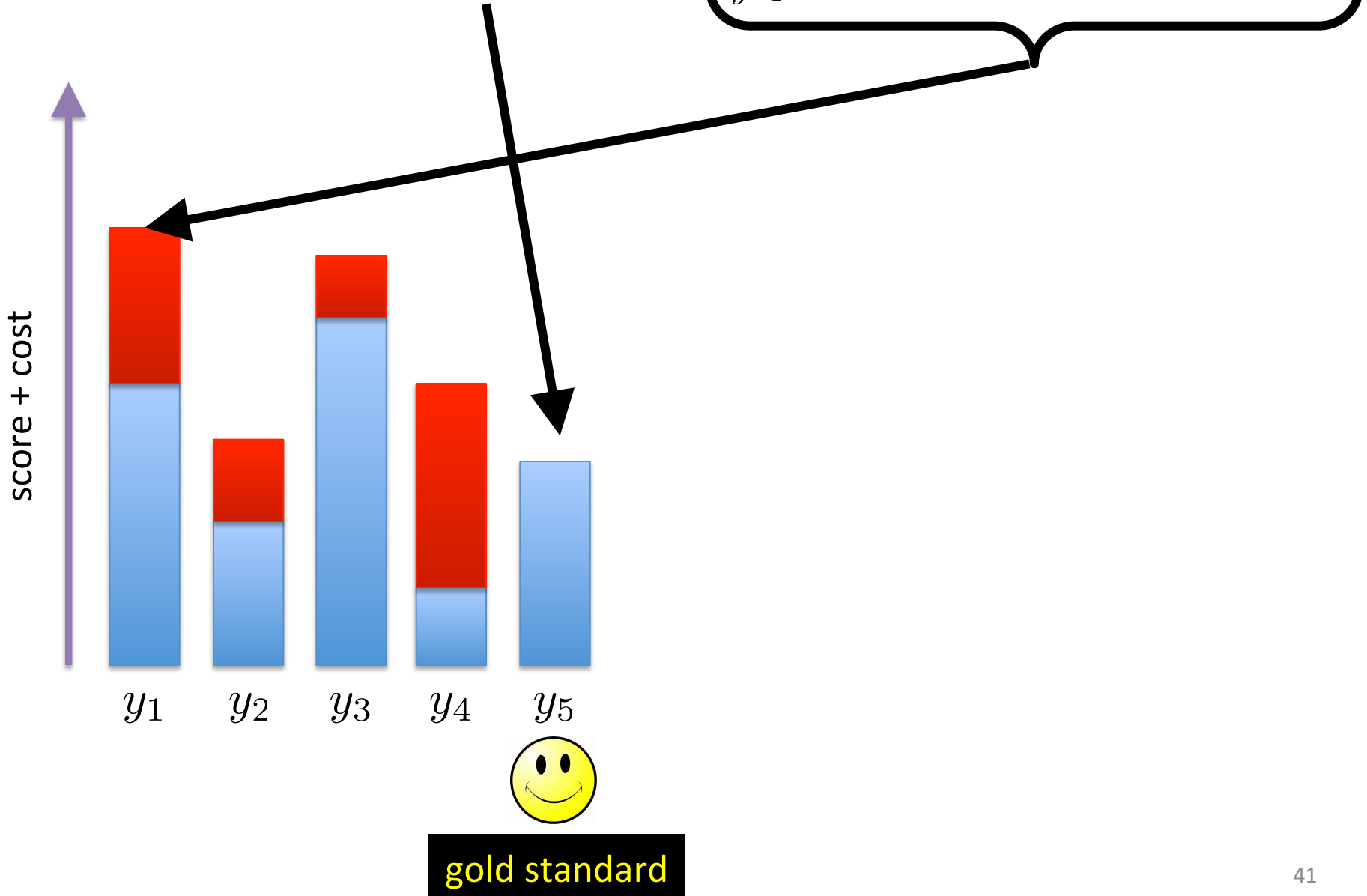
hinge loss:

$$\text{loss}_{\text{hinge}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w}) + \underbrace{\max_{y' \in \mathcal{L}} (\text{score}(\mathbf{x}, y', \mathbf{w}) + \text{cost}(y, y'))}_{\text{max over } y' \in \mathcal{L}}$$



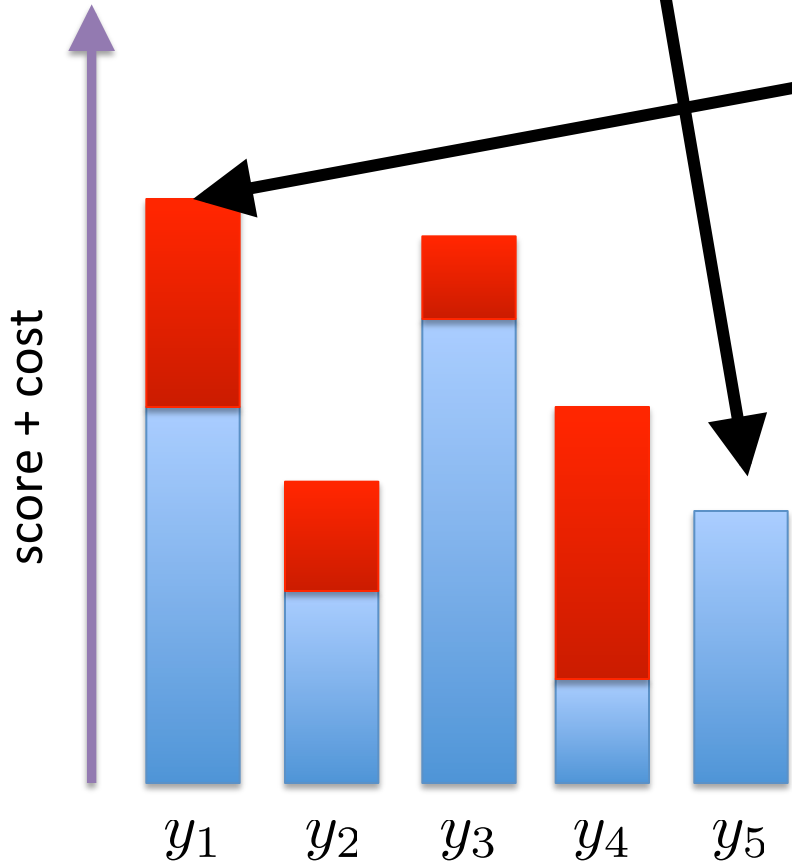
hinge loss:

$$\text{loss}_{\text{hinge}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w}) + \max_{y' \in \mathcal{L}} (\text{score}(\mathbf{x}, y', \mathbf{w}) + \text{cost}(y, y'))$$



hinge loss:

$$\text{loss}_{\text{hinge}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w}) + \max_{y' \in \mathcal{L}} (\text{score}(\mathbf{x}, y', \mathbf{w}) + \text{cost}(y, y'))$$



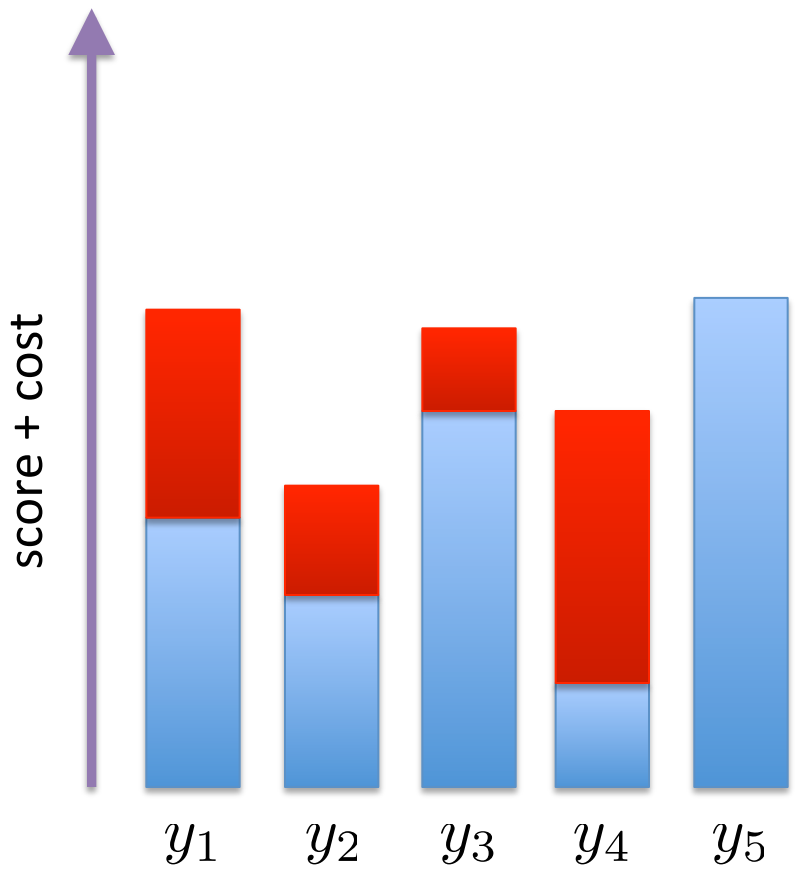
effect of learning?



gold standard

hinge loss:

$$\text{loss}_{\text{hinge}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w}) + \max_{y' \in \mathcal{L}} (\text{score}(\mathbf{x}, y', \mathbf{w}) + \text{cost}(y, y'))$$

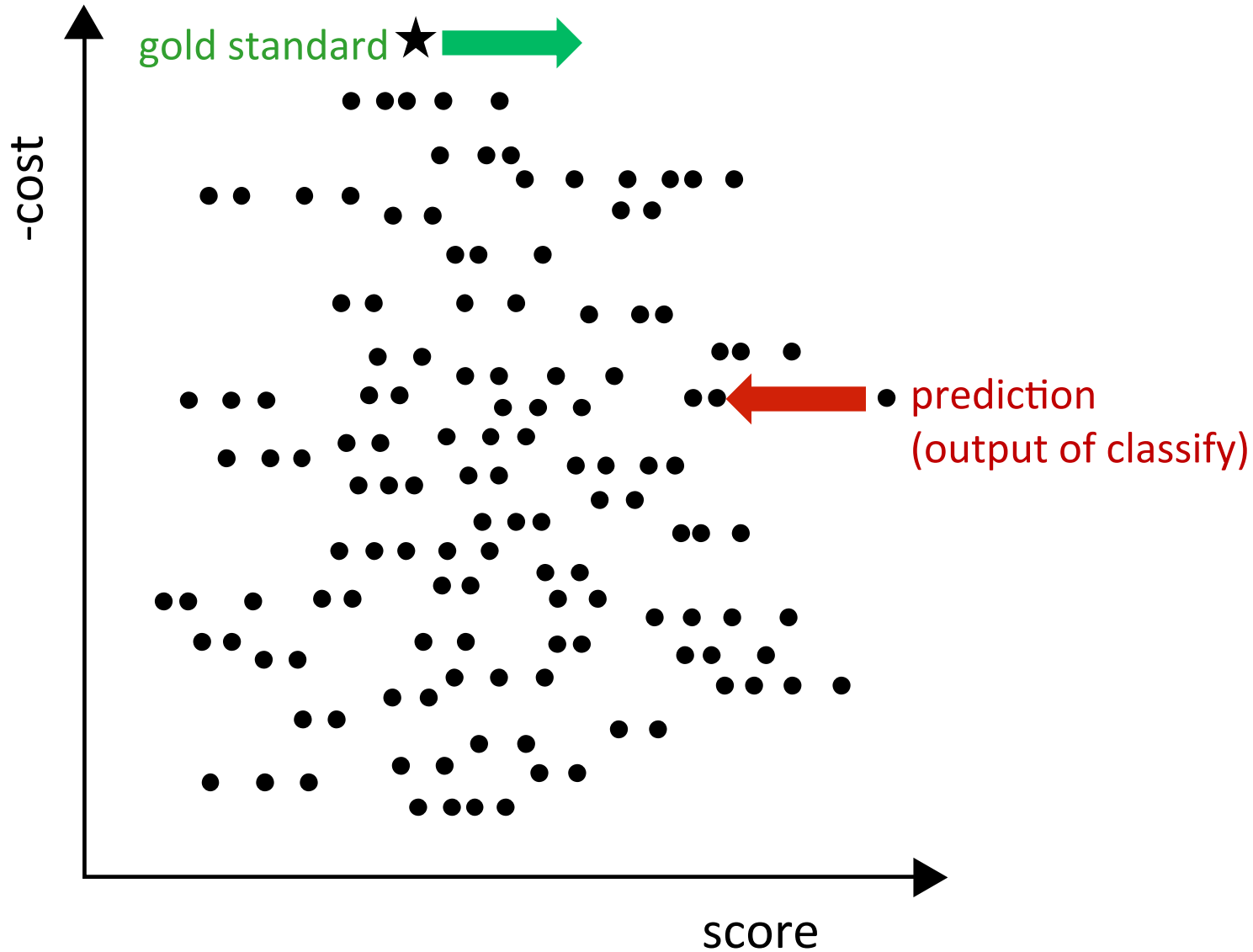


effect of learning:
score of gold standard
will be higher than
score+cost of all
others

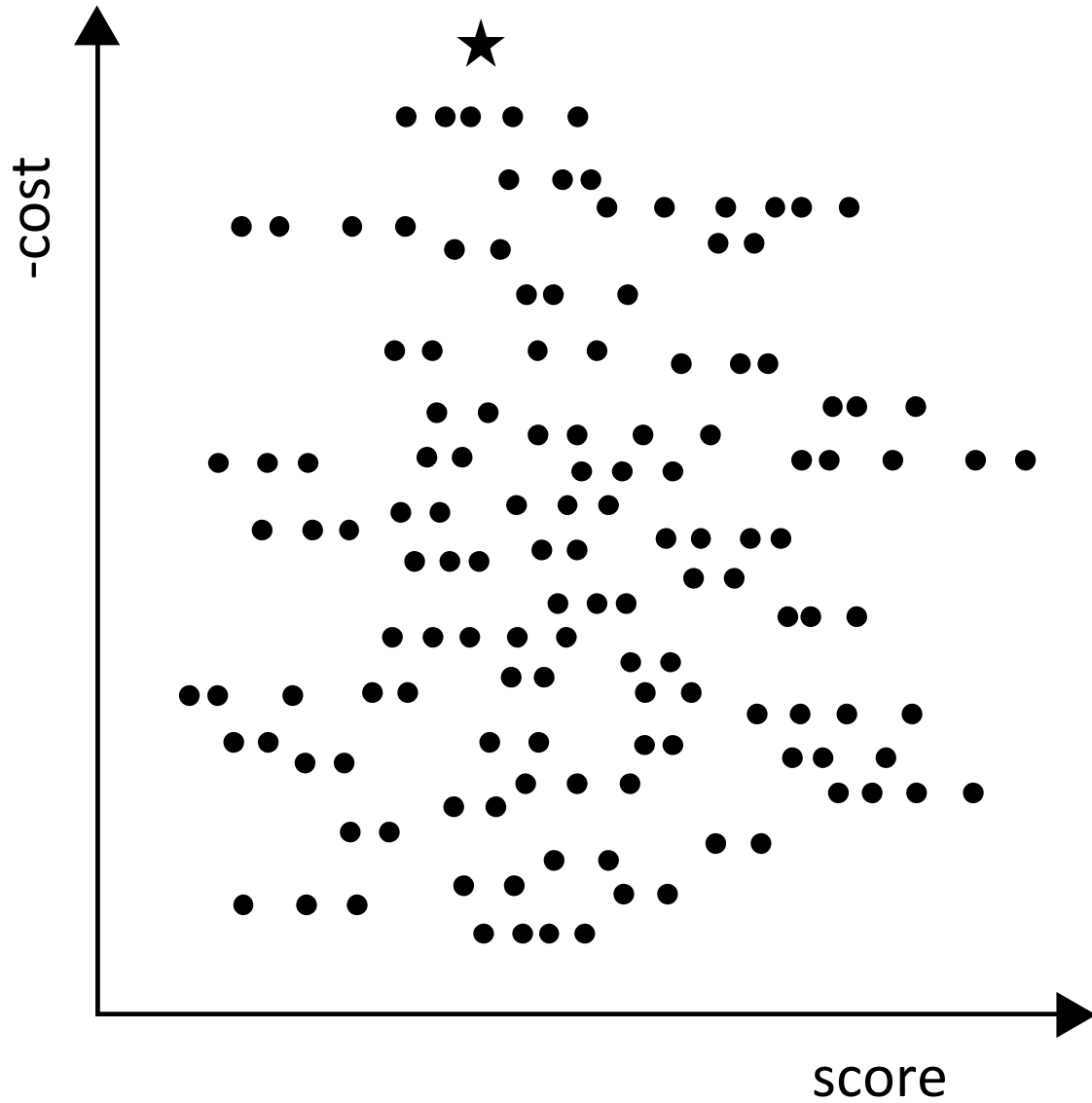


gold standard

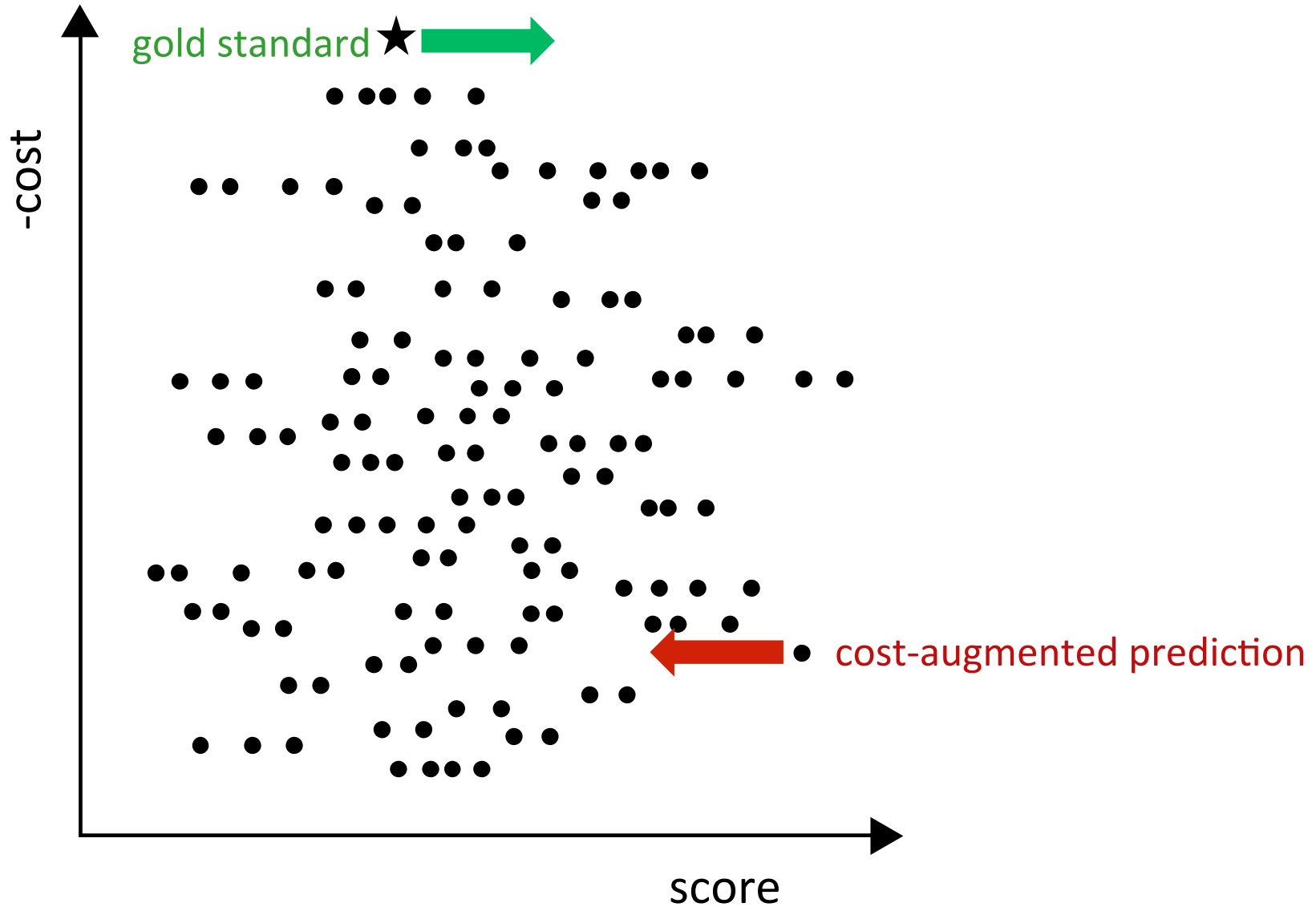
Perceptron Loss



Hinge Loss?



Hinge Loss



Perception \rightarrow Hinge

$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = - \sum_i w_i f_i(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \sum_i w_i f_i(\mathbf{x}, y')$$

$$\text{loss}_{\text{hinge}}(\mathbf{x}, y, \mathbf{w}) = - \sum_i w_i f_i(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \left(\sum_i w_i f_i(\mathbf{x}, y') + \text{cost}(y, y') \right)$$

Loss Subgradients for Linear Models

- hinge loss:

$$\text{loss}_{\text{hinge}}(\mathbf{x}, y, \mathbf{w}) = - \sum_i w_i f_i(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \left(\sum_i w_i f_i(\mathbf{x}, y') + \text{cost}(y, y') \right)$$

- subderivative for a single parameter:

Loss Subgradients for Linear Models

- hinge loss:

$$\text{loss}_{\text{hinge}}(\mathbf{x}, y, \mathbf{w}) = - \sum_i w_i f_i(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \left(\sum_i w_i f_i(\mathbf{x}, y') + \text{cost}(y, y') \right)$$

- subderivative for a single parameter:

$$\frac{\partial \text{loss}_{\text{hinge}}(\mathbf{x}, y, \mathbf{w})}{\partial w_j} = -f_j(\mathbf{x}, y) + f_j(\mathbf{x}, \text{costClassify}(\mathbf{x}, y, \mathbf{w}))$$

Loss Subgradients for Linear Models

- hinge loss:

$$\text{loss}_{\text{hinge}}(\mathbf{x}, y, \mathbf{w}) = - \sum_i w_i f_i(\mathbf{x}, y) + \max_{y' \in \mathcal{L}} \left(\sum_i w_i f_i(\mathbf{x}, y') + \text{cost}(y, y') \right)$$

- subderivative for a single parameter:

$$\frac{\partial \text{loss}_{\text{hinge}}(\mathbf{x}, y, \mathbf{w})}{\partial w_j} = -f_j(\mathbf{x}, y) + f_j(\mathbf{x}, \text{costClassify}(\mathbf{x}, y, \mathbf{w}))$$

$$\text{costClassify}(\mathbf{x}, y, \mathbf{w}) = \operatorname{argmax}_{y' \in \mathcal{L}} \left(\sum_i w_i f_i(\mathbf{x}, y') + \text{cost}(y, y') \right)$$

“cost-augmented inference” or “cost-augmented decoding”

Feature count cut-off of zero?

- perceptron loss update rule:

$$w_j \leftarrow w_j + \eta (f_j(\mathbf{x}, y) - f_j(\mathbf{x}, \text{classify}(\mathbf{x}, \mathbf{w})))$$

- what do you expect to happen to weights of features with count 0 in the training data? (if they are initialized to 0)

Feature count cut-off of zero?

- perceptron loss update rule:

$$w_j \leftarrow w_j + \eta (f_j(\mathbf{x}, y) - f_j(\mathbf{x}, \text{classify}(\mathbf{x}, \mathbf{w})))$$

- what do you expect to happen to weights of features with count 0 in the training data? (if they are initialized to 0)
 - they will stay at zero or become negative

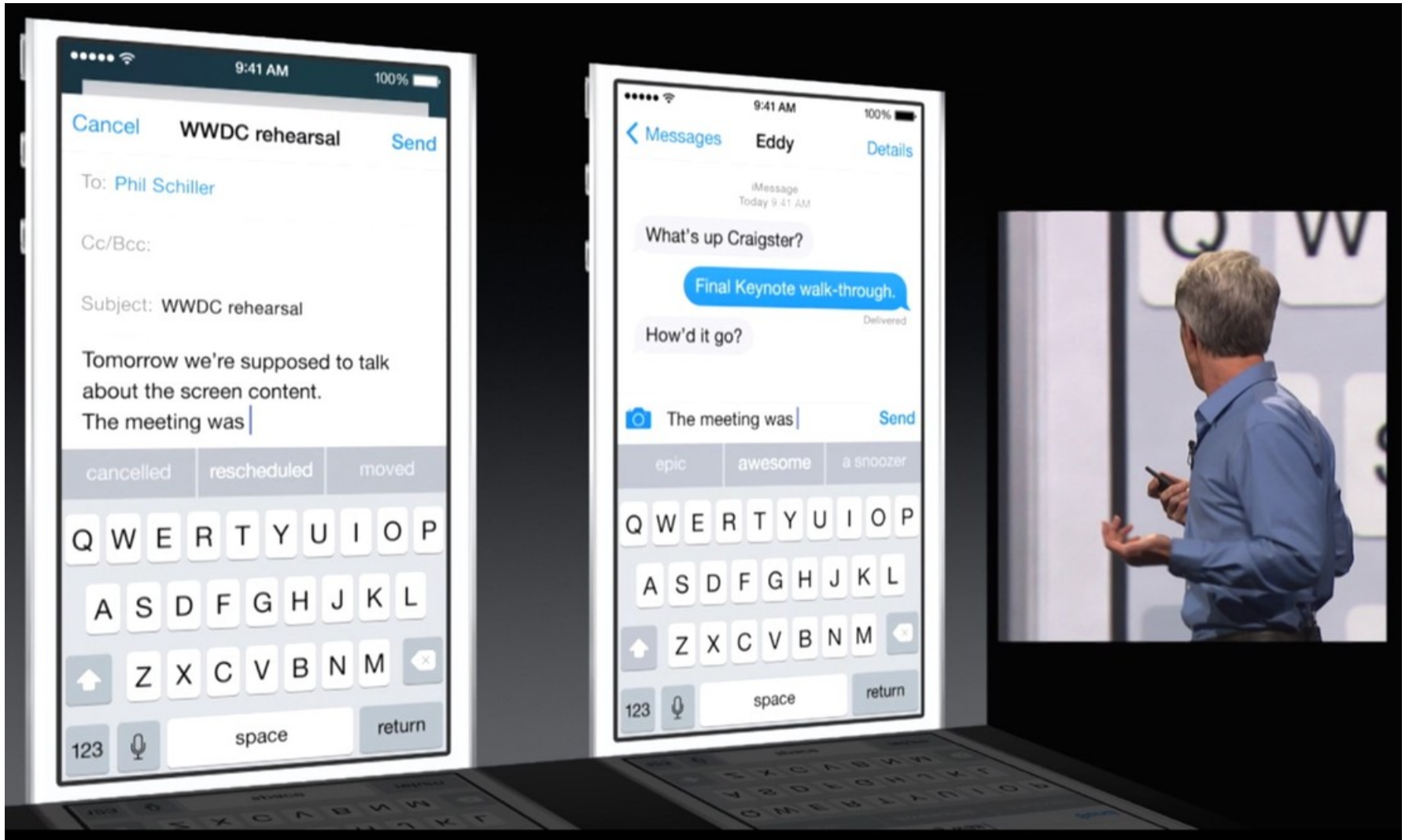
Roadmap

- words, morphology, lexical semantics
- text classification
- simple neural methods for NLP
- language modeling and word embeddings
- recurrent/recursive/convolutional networks in NLP
- sequence labeling, HMMs, dynamic programming
- syntax and syntactic parsing
- semantics, compositionality, semantic parsing
- machine translation and other NLP tasks

Probabilistic Language Models

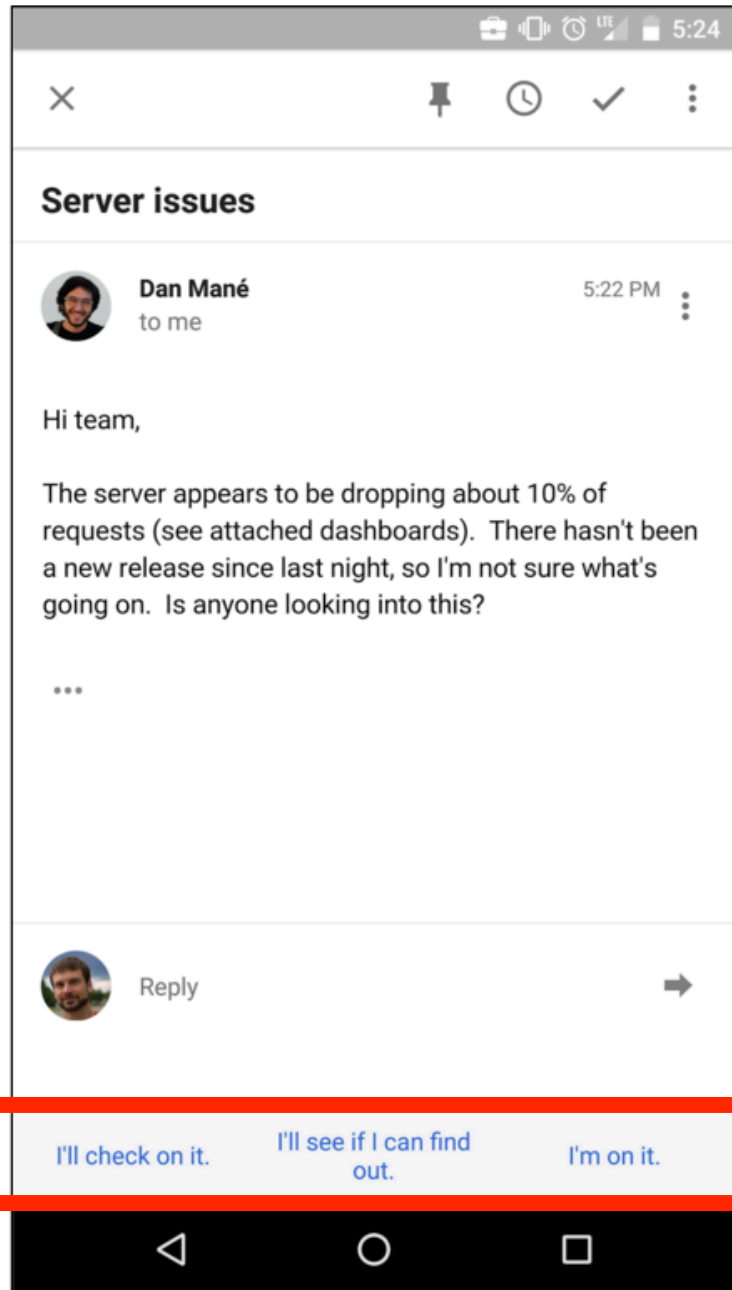
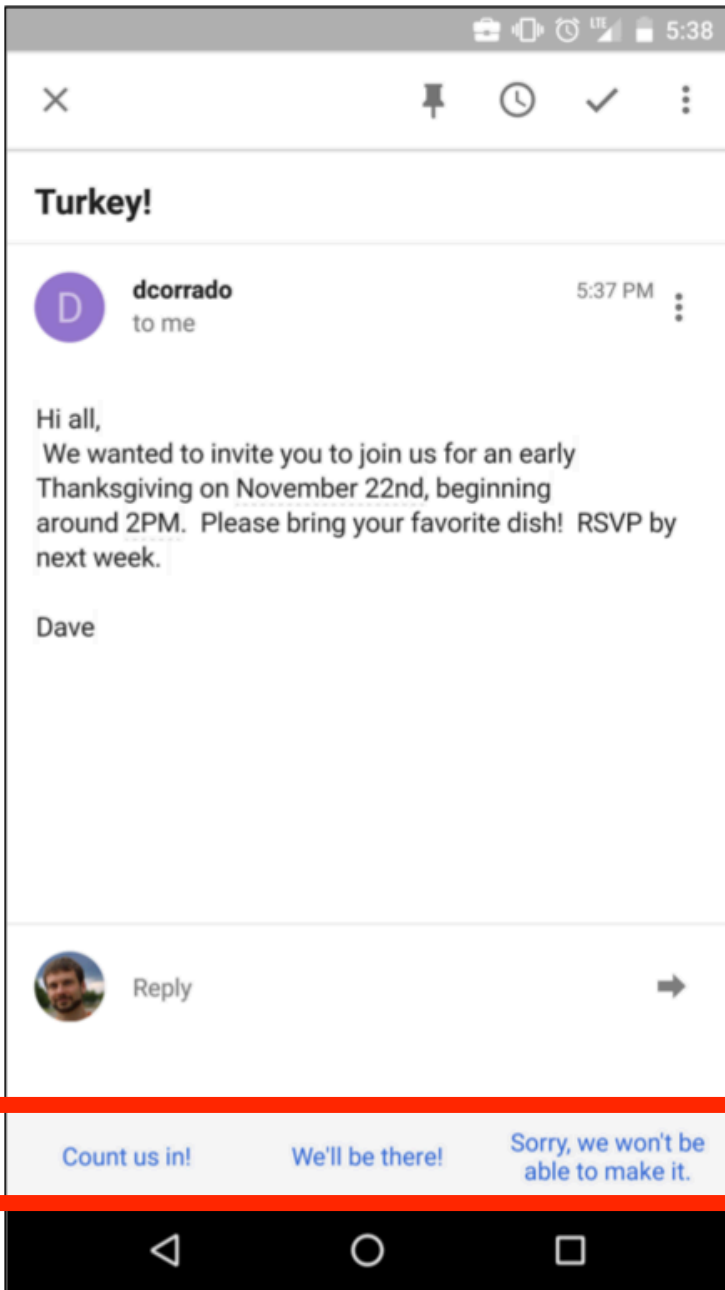
- language modeling: assign probabilities to sentences
- Why?
 - machine translation:
 - $P(\text{high winds tonite}) > P(\text{large winds tonite})$
 - spelling correction:
 - The office is about fifteen **minuets** from my house
 - $P(\text{about fifteen minutes from}) > P(\text{about fifteen minuets from})$
 - speech recognition:
 - $P(\text{I saw a van}) \gg P(\text{eyes awe of an})$
 - summarization, question answering, etc.!

Automatic Completion



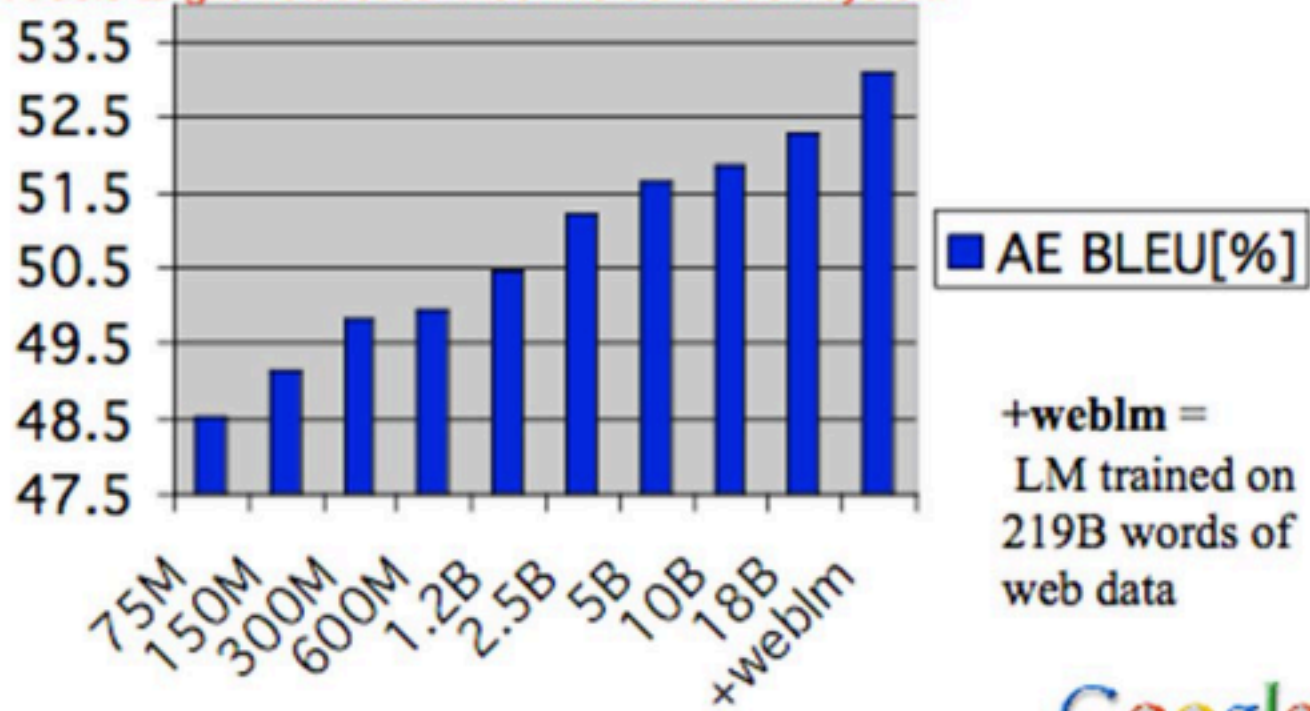
Automatic Completion





Language Modeling for Machine Translation

Impact on size of language model training data (in words) on quality of Arabic-English statistical machine translation system



Probabilistic Language Modeling

- goal: compute the probability of a sequence of words:

$$P(\mathbf{w}) = P(w_1, w_2, \dots, w_n)$$

- related task: probability of next word:

$$P(w_4 \mid w_1, w_2, w_3)$$

- a model that computes either of these:

$$P(\mathbf{w}) \quad \text{or} \quad P(w_k \mid w_1, w_2, \dots, w_{k-1})$$

is called a **language model (LM)**

How to compute $P(w)$

- How to compute this joint probability:
 - $P(\text{its, water, is, so, transparent, that})$
- Intuition: let's rely on the Chain Rule of Probability

Reminder: Chain Rule

- factor joint probability into product of conditional probabilities:

$$P(w_1, w_2, \dots, w_n) = P(w_1)P(w_2 | w_1)P(w_3 | w_1, w_2) \dots P(w_n | w_1, w_2, \dots, w_{n-1})$$

- we have not yet made any independence assumptions

Chain Rule for computing joint probability of words in sentence

$$P(w_1, w_2, \dots, w_n) = \prod_i P(w_i \mid w_1, w_2, \dots, w_{i-1})$$

$P(\textit{“its water is so transparent”}) =$

$P(\textit{its}) \times P(\textit{water} \mid \textit{its}) \times P(\textit{is} \mid \textit{its water})$

$\times P(\textit{so} \mid \textit{its water is}) \times P(\textit{transparent} \mid \textit{its water is so})$

How to estimate these probabilities

- could we just count and divide?

$$P(\text{the l its water is so transparent that}) = \frac{\textit{Count}(\text{its water is so transparent that the})}{\textit{Count}(\text{its water is so transparent that})}$$

- no! too many possible sentences!
- we'll never see enough data for estimating these

Markov Assumption



Andrei Markov

- simplifying assumption:

$P(\text{the } l \text{ its water is so transparent that}) \approx P(\text{the } l \text{ that})$

- or maybe:

$P(\text{the } l \text{ its water is so transparent that}) \approx P(\text{the } l \text{ transparent that})$

Markov Assumption

- i.e., we approximate each component in the product:

$$P(w_i \mid w_1, \dots, w_{i-2}, w_{i-1}) \approx P(w_i \mid w_{i-k}, \dots, w_{i-2}, w_{i-1})$$

Simplest case: Unigram model

$$P(w_1, w_2, \dots, w_n) = \prod_i P(w_i)$$

automatically generated sentences from a unigram model:

fifth an of futures the an incorporated a a the
inflation most dollars quarter in is mass

thrift did eighty said hard 'm july bullish

that or limited the

Bigram model

condition on the previous word:

$$P(w_1, w_2, \dots, w_n) = \prod_i P(w_i \mid w_{i-1})$$

automatically generated sentences from a bigram model:

texaco rose one in this issue is pursuing growth in a boiler
house said mr. gurria mexico 's motion control proposal
without permission from five hundred fifty five yen

outside new car parking lot of the agreement reached

this would be a record november

n-gram models

- we can extend to trigrams, 4-grams, 5-grams
- in general this is an insufficient model of language
 - because language has **long-distance dependencies**:

“The computer which I had just put into the machine room on the fifth floor crashed.”

- but we can often get away with n-gram models