# Large-scale machine learning
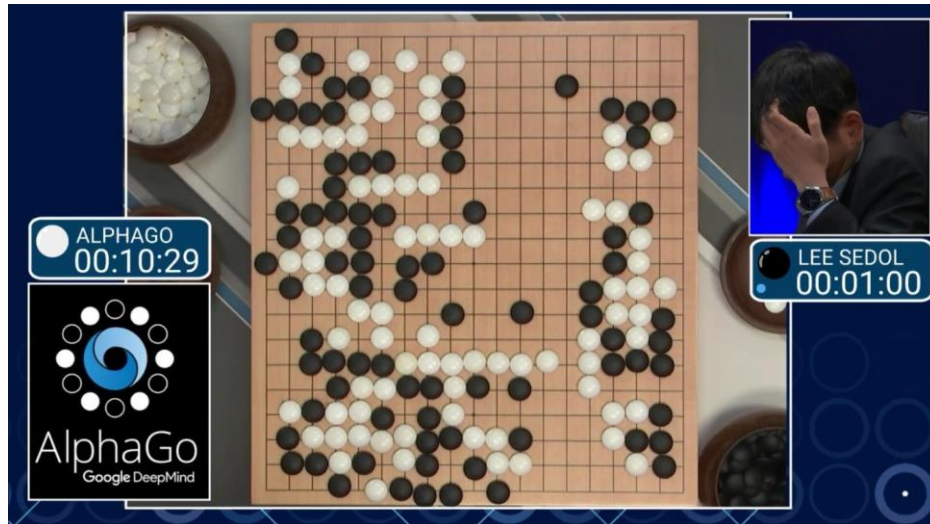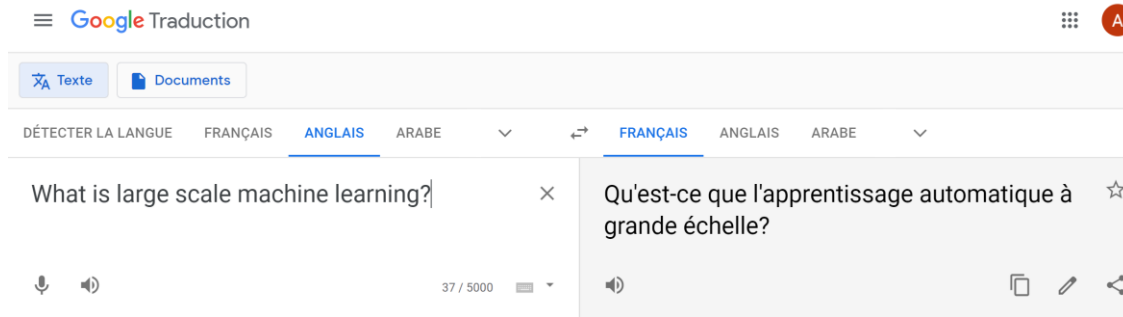
**D**    what is large scale machine learning?

Large scale machine learning refers to the process of training and deploying machine learning models on very large datasets. In this context, "large" typically means datasets that are too big to fit into the memory of a single machine, and which may require distributed computing resources to process.

Large scale machine learning is important for many applications in industry, science, and engineering, where the volume of data is growing rapidly and the complexity of models is increasing. Examples include natural language processing, computer vision, recommendation systems, and fraud detection.

To deal with the challenges of large scale machine learning, researchers and practitioners have developed a range of algorithms, tools, and platforms that enable efficient processing and training of models on distributed systems. These include frameworks like Apache Spark, TensorFlow, and PyTorch, as well as specialized hardware like graphics processing units (GPUs) and tensor processing units (TPUs).

# Recent ML breakthroughs

**Course overview**

Session 1&2: Introduction to large-scale machine learning
Session 3&4: Distributed optimization
Session 5&6: Distributed representation learning
Session 7&8: Project/evaluation, *To be defined!*

Each session will combine theory and practice.

# Teachers



**David Diebold**
*Software Development Engineer*

**Etienne Duchesne**
*PHD, Machine Learning Engineer*

**Martin Bompaire**
*PHD, Machine Learning Engineer*

**Alexandre Gilotte**
*PHD, Researcher*

# Acknowledgement

Those slides are largely inspired / copied from several sources:

Léon Bottou's "Large-scale machine learning revisited" conference:
https://bigdata2013.sciencesconf.org/conference/bigdata2013/pages/bottou.pdf

Sanjiv Kumar's "Large-scale machine learning" course:
http://www.sanjivk.com/EECS6898/lectures.html

Jean-Philippe Vert's "Large-Scale Machine Learning" course:
http://members.cbio.mines-paristech.fr/~jvert/svn/lsml/lsml18/

# Content


Tradeoffs


More Data ?


Labelling


Scaling Tricks

# Tradeoffs

# Example: Linear regression

n examples

$X$ * $w$ = $y$

d features

$\hat{w}$ minimizing quadratic error:

$$L(w) = \sum (x_i.w - y_i)^2 + \lambda w^T w$$
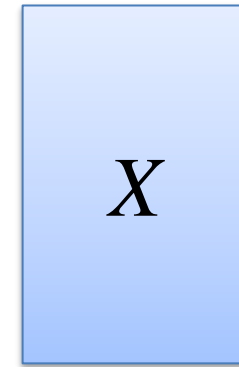$$= (Xw - y)^T (Xw - y) + \lambda w^T w$$

Solution: $\hat{w} = (X^T X + \lambda I)^{-1} X^T y$

# Example: Linear regression

$$\hat{w} = (X^T X + \lambda I)^{-1} X^T y$$

n examples

$X$

d features

```
1  XXt = X.transpose().dot( X ) + l2_regularization
2  w =  np.linalg.inv(XXt) .dot(X.transpose() ).dot(  Y )
```
executed in 1ms, finished 18:40:31 2021-03-02

```
1  predictions = X.dot(w)
```
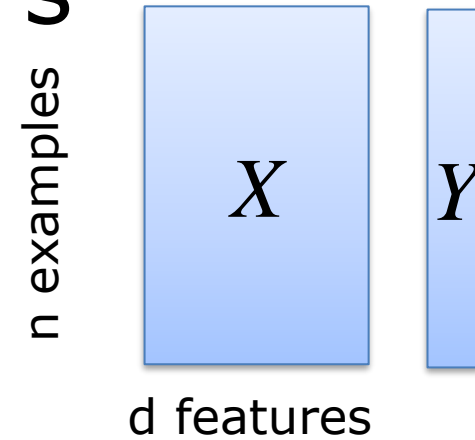executed in 3ms, finished 18:40:31 2021-03-02

O(d³)

O(n.d²)

What if n = $10^7$
and d = $10^6$ ?

# Let's generalize
# Supervised learning reminders

- Data: independent examples $(X_i, Y_i)$
- Goal: find f such that $Y \approx f(X)$ ?

- Formally, looking for f minimizing average « loss »

$$f^* := \underset{f}{\text{Argmin}} \left( \mathbb{E}( \text{loss}( f(X), Y ) ) \right)$$

n examples

$X$    $Y$

d features

On average, when X and Y follow the unknown distribution of the datset

Loss measuring the "error" between prediction and label. Example: mean square error, loglikelihood ,…

# Supervised learning reminders

$$f^* := \text{Argmin} \, (\mathbb{E}( \text{loss}( f(X), Y \, ) \, )$$

f ∈ All functions

Not realistic to find minimizer among all possible functions!

Cannot compute! The true distribution is unknown.

n examples

$X$   $Y$

d features

- Instead, minimize in a class $\mathcal{F}$ of parametric functions

- $f^*_{\mathcal{F}} := \text{Argmin} \, (\mathbb{E}( \text{loss}( f(X), Y))$

  f ∈ $\mathcal{F}$

- Approximation error:

  $\mathbb{E}(\text{loss}(f^*_{\mathcal{F}})) - \mathbb{E}(\text{loss}(f^*))$

  Because $f^*$ not in $\mathcal{F}$

- Instead, approximate by average on training set:

  $\mathbb{E}( \text{loss}( f(X), Y \, ) \, ) \approx 1/n \, \Sigma_i \, \text{loss}( f(X_i), Y_i)$

- $f_n := \text{Argmin} \, (1/n \, \Sigma_i \, \text{loss}( f(X_i), Y_i) \, )$

  f ∈ $\mathcal{F}$

- Estimation error:

  $\mathbb{E}( \text{loss}(f_n)) - \mathbb{E}(\text{loss}(f^*_{\mathcal{F}}) \, )$

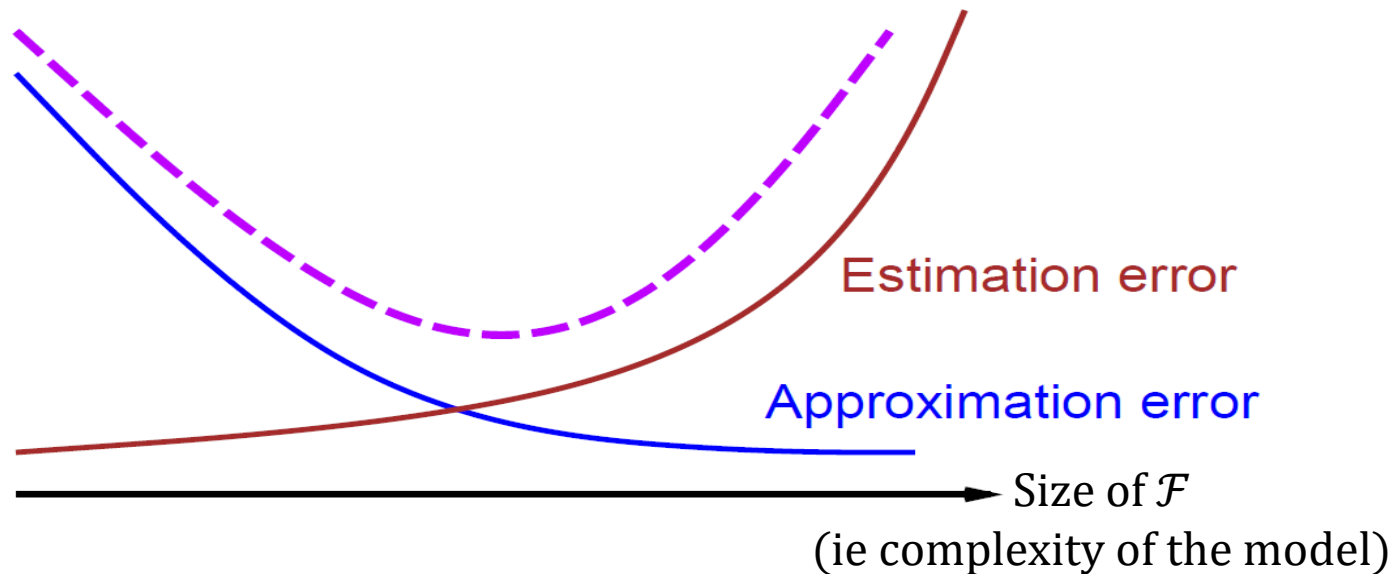  Not enough data to identify $f^*_{\mathcal{F}}$

# Model selection tradeoffs

$$f^* := \text{Argmin} \left( \mathbb{E}( \text{loss}( f(X), Y ) ) \right)$$

**For a given dataset of size n** ——————→ $f_n := \text{Argmin} \left( 1/n \sum_i \text{loss}( f(X_i), Y_i ) \right)$
$$f \in \mathcal{F}$$

Error decomposition:

$| \mathbb{E}(\text{loss}( f_n )) - \mathbb{E}(\text{loss}( f^* )) |$ = $| \mathbb{E}(\text{loss}( f^*_{\mathcal{F}} )) - \mathbb{E}(\text{loss}(f^* )) |$   Approximation error

$+ | \mathbb{E}(\text{loss}(f_n )) - \mathbb{E}(\text{loss}( f^*_{\mathcal{F}} )) |$   Estimation error



Estimation error

Approximation error

Size of $\mathcal{F}$
(ie complexity of the model)

How complex a model can you afford with your data?

# Learning with approximate optimization

Optimization problem.
Costly to solve accurately.

$$f_n := \text{Argmin}_{f \in \mathcal{F}} \left( 1/n \sum_i \text{loss}( f(X_i), Y_i) \right)$$

- Instead: define stopping criteria $\rho$
- Let $\hat{f}_n$ the approximate solution returned by the optimizer.

- Error decomposition:

$$| \mathbb{E}(\text{loss}( \hat{f}_n )) - \mathbb{E}(\text{loss}(f^*)) | = | \mathbb{E}(\text{loss}(f^*_{\mathcal{F}} )) - \mathbb{E}(\text{loss}(f^*)) | \quad \text{Approximation error}$$
$$+ | \mathbb{E}(\text{loss}(f_n )) - \mathbb{E}(\text{loss}(f^*_{\mathcal{F}})) | \quad \text{Estimation error}$$
$$+ | \mathbb{E}(\text{loss}(\hat{f}_n )) - \mathbb{E}(\text{loss}(f_n)) | \quad \text{Optimization error}$$

- Choose $\mathcal{F}$, n, $\rho$ to get small total error
- Subject to constraints: number of available samples, max compute time.

# Small scale versus large scale

## Small scale learning problem

- We have a small-scale learning problem when the active budget constraint is the number of examples $n$.

## Large-scale learning problem

- We have a large-scale learning problem when the active budget constraint is the computing time $T$.
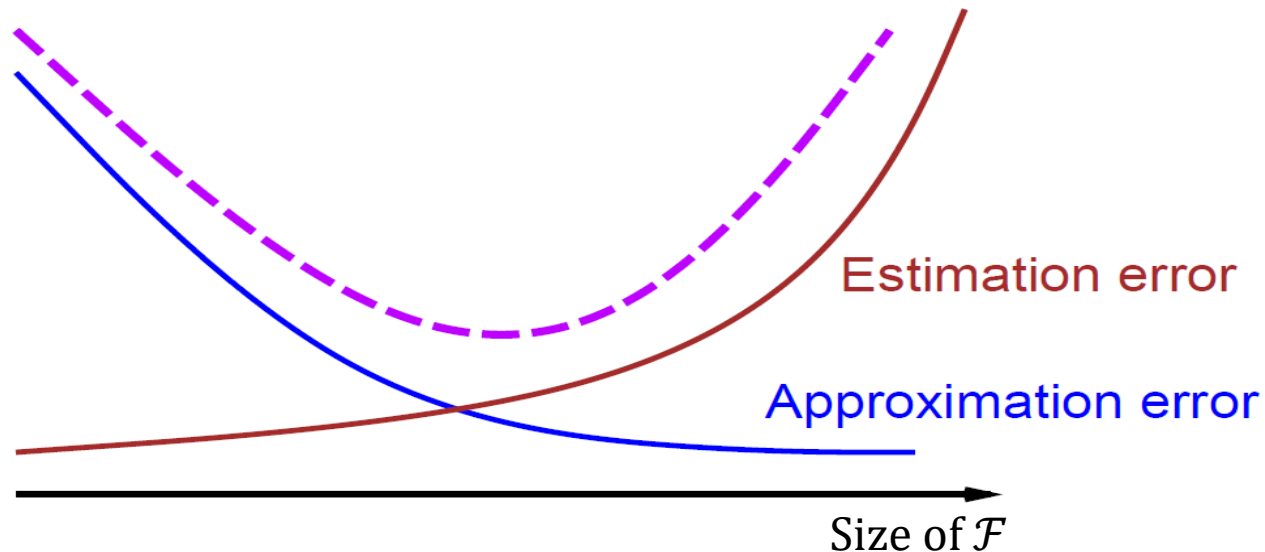
# Small scale learning

Constrained by the number of samples
➔ Use all samples to reduce estimation error

Two degrees of freedom:
- Optimize as precisely as possible to avoid significative optimization error
- Select the size of $\mathcal{F}$ by crossvalidation



Estimation error

Approximation error

Size of $\mathcal{F}$

# Large scale learning

Constrained by the training time / training ressources
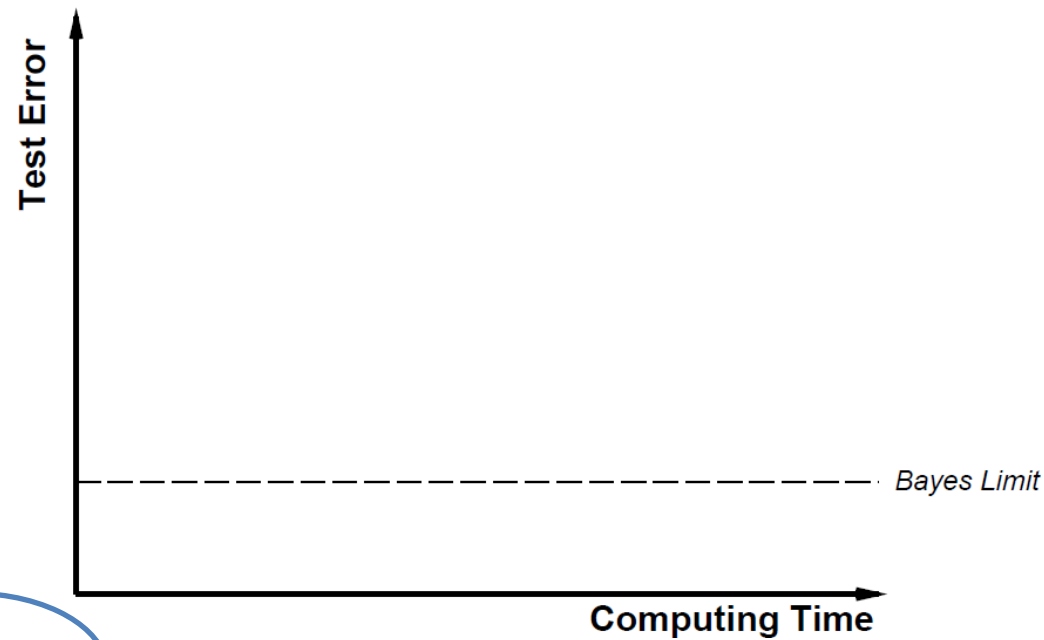
Example:  OpenAI GPT-3 text model

- Training set: 45 TB text data, mostly crawled from internet.
- Training cost estimated to several M€

Computing time depends on n, $\mathcal{F}$ and ρ.

- Should you use more samples or spend more time optimizing on smaller sample set?
- Methods to reduce dataset size with small loss of precision may *improve* final performances!
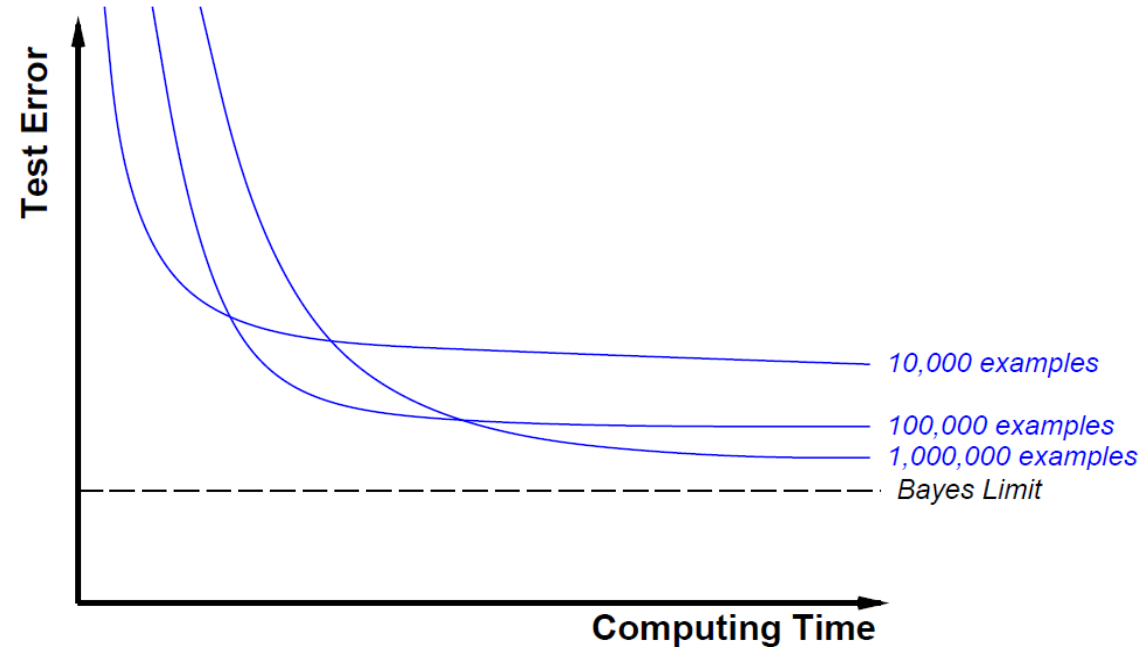- Best tradeof also depends on optimization algorithm.

# Test error versus computing time



Test error:

- $\frac{1}{n_{test}} \sum_{i \in TestSet} loss(\hat{f}_n(X_i), Y_i)$

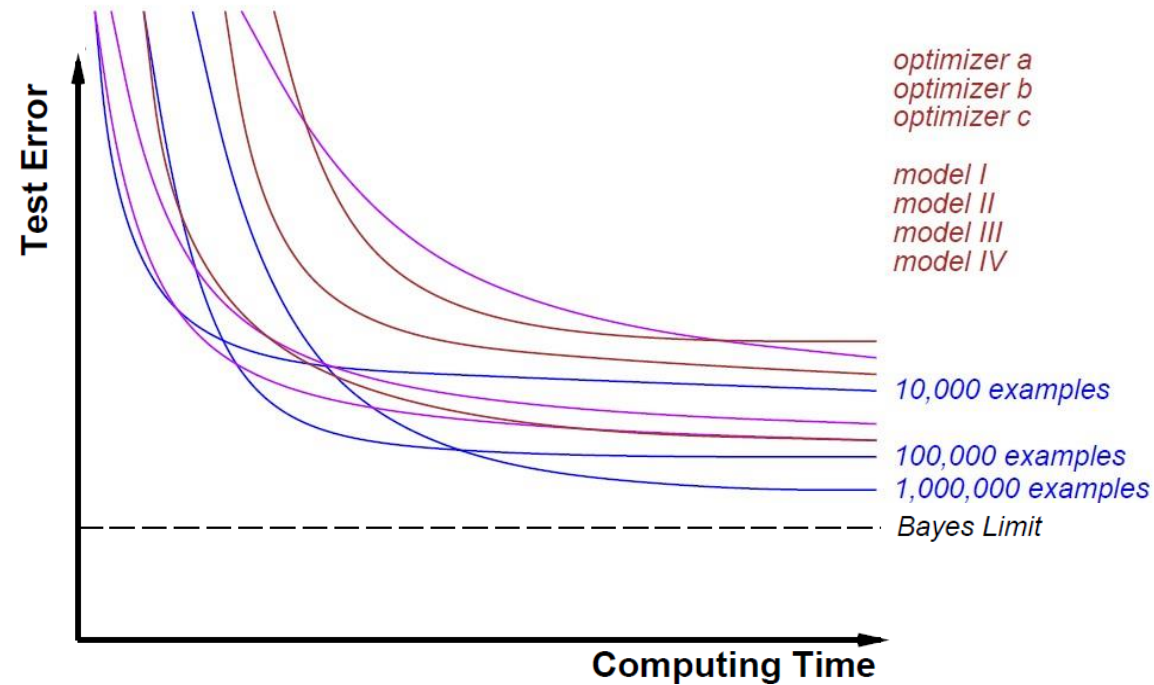- Unbiased estimator of $\mathbb{E}(loss(\hat{f}_n))$

# Test error versus computing time



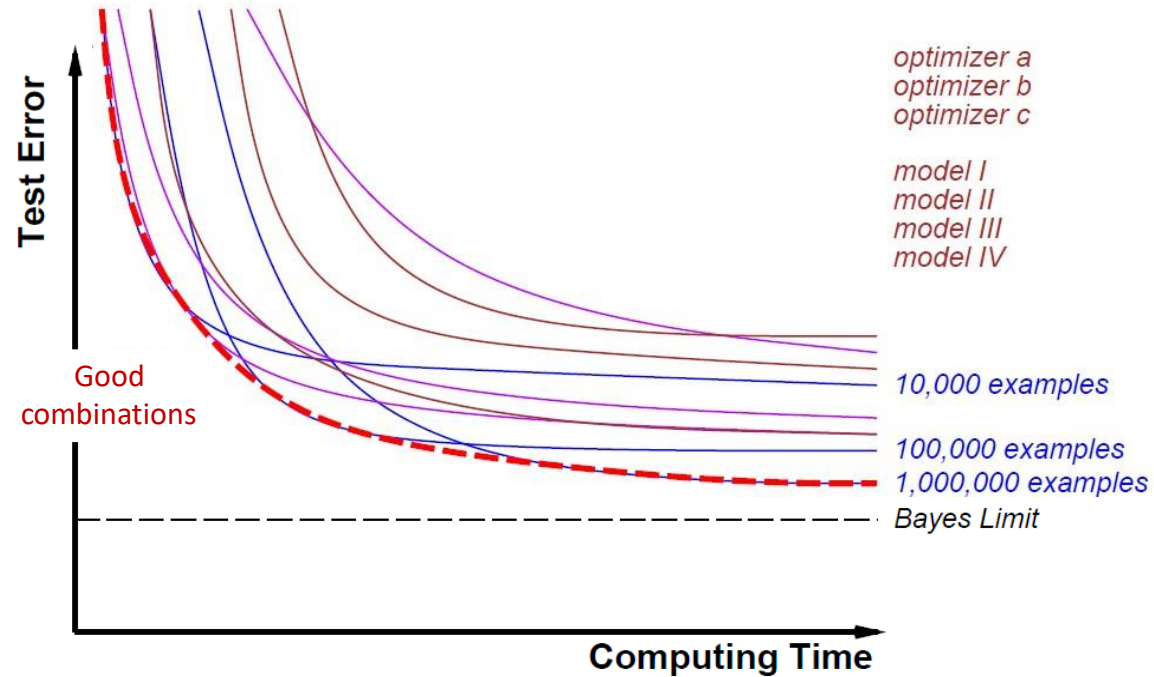Pushing optimization further, not changing n or $\mathcal{F}$

- Vary the number of examples

# Test error versus computing time



optimizer a
optimizer b
optimizer c

model I
model II
model III
model IV

10,000 examples

100,000 examples
1,000,000 examples
Bayes Limit

Test Error

Computing Time

- Vary the number of examples, the model, the algorithm

# Test error versus computing time



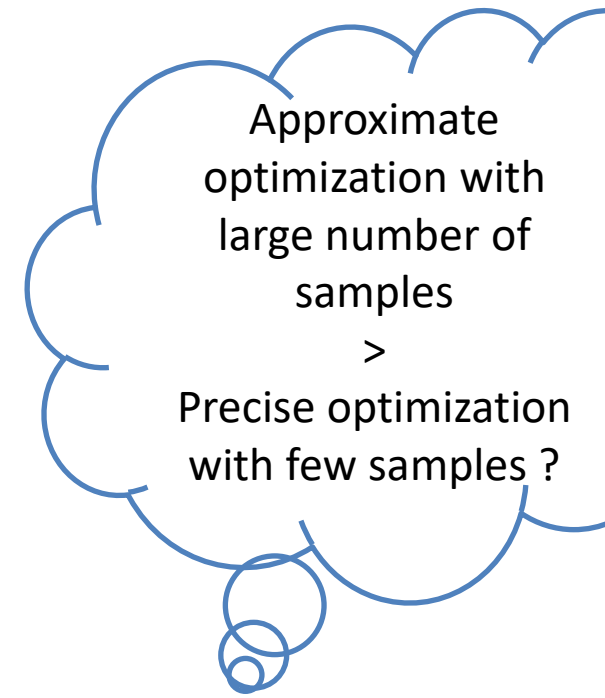- Optimal combination depends on training time budget.

# The tradeoffs of large-scale learning

Small-scale learning ≠ large-scale learning

- Large-scale learning involves more complex tradeoffs that depends on the properties of the optimization algorithm.

Good optimization algorithm ≠ good learning algorithm

- Mediocre optimization algorithms (e.g., SGD) often outperform sophisticated optimization algorithms on large-scale learning problems.
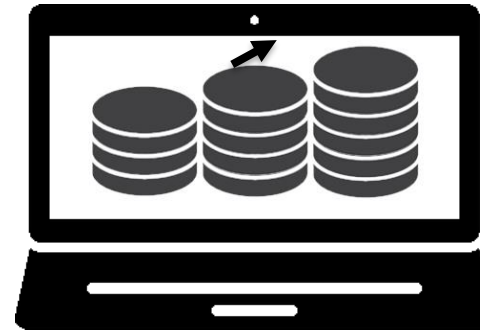
Approximate optimization with large number of samples
>
Precise optimization with few samples ?

# Content



Tradeoffs

More Data ?

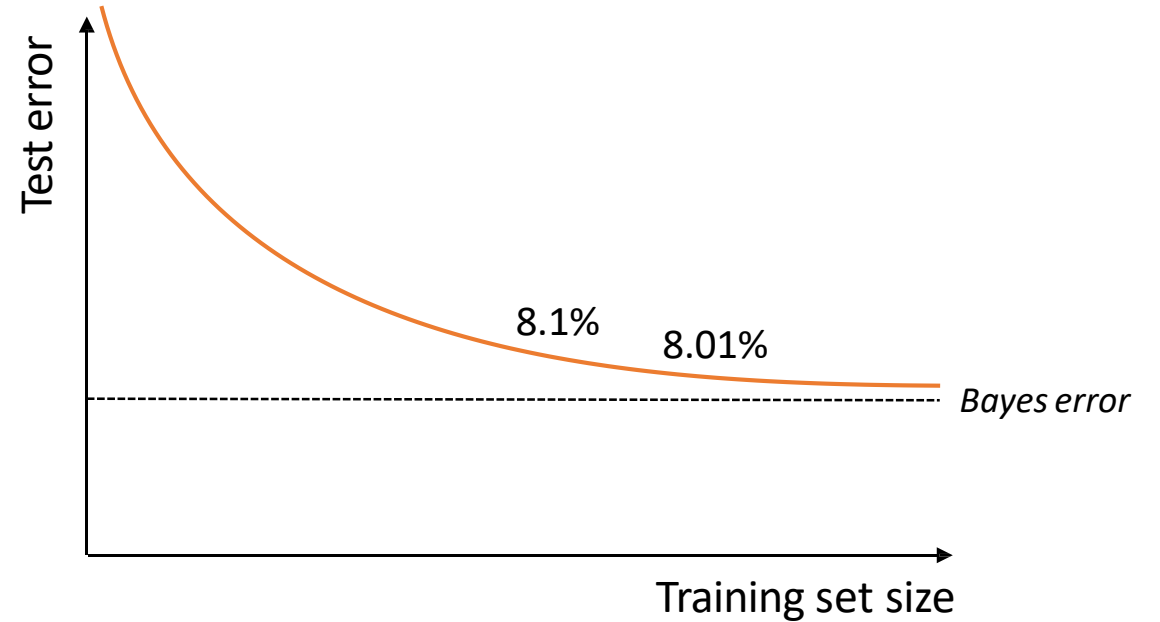Labelling

Scaling Tricks

# More Data ?

# Context

**Dataset size:** n
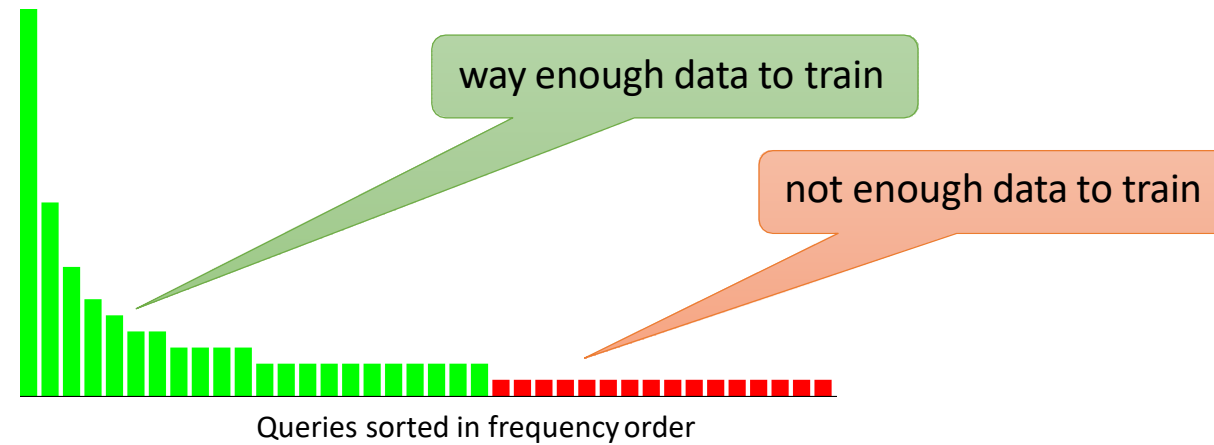**Training time**: 12 hours (a bit too much but ok)
**Loss**: 8.01%



**Question**
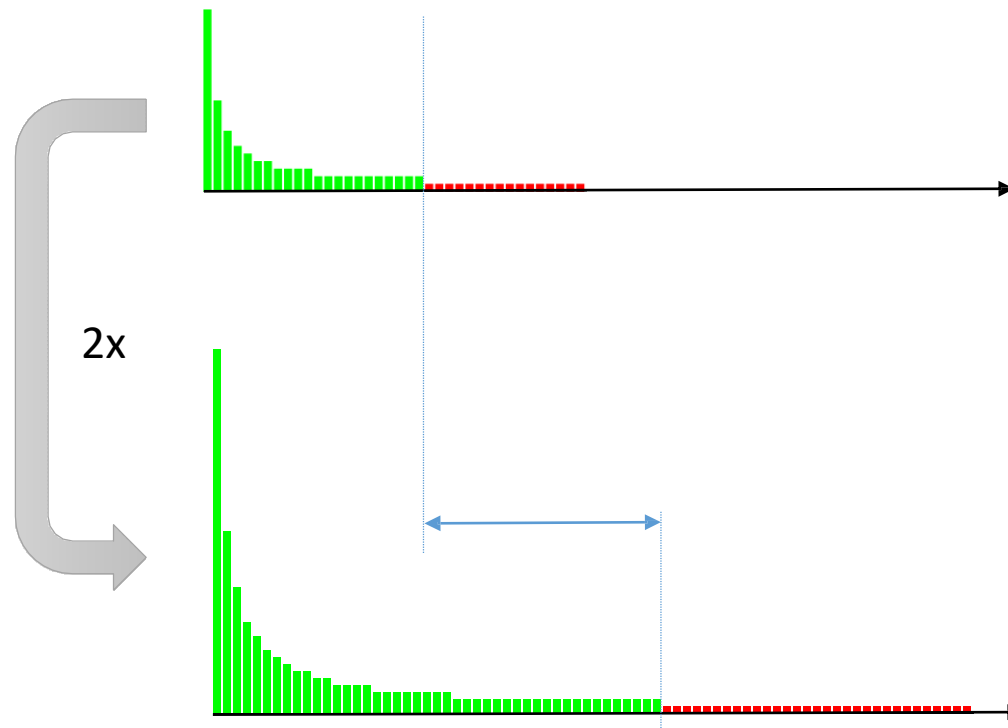Opportunity to use a dataset twice bigger.
What would you do ?

# Zipf distributed data

Roughly half of the search queries are unique.



Queries sorted in frequency order

# Doubling the size of the training set



**Average error** not much improved: Most examples in the test set are from the head of the distribution, and already well predicted.

**Proportion of distinct querries** we learn to answer correctly increased a lot!

# Value of big data

Accuracy improvements are subject to diminishing returns.

Breadth improvements are not subject to diminishing returns.

*"How accurately do we recognize an object category?"*
      *vs. "How many categories do we recognize well enough?"*

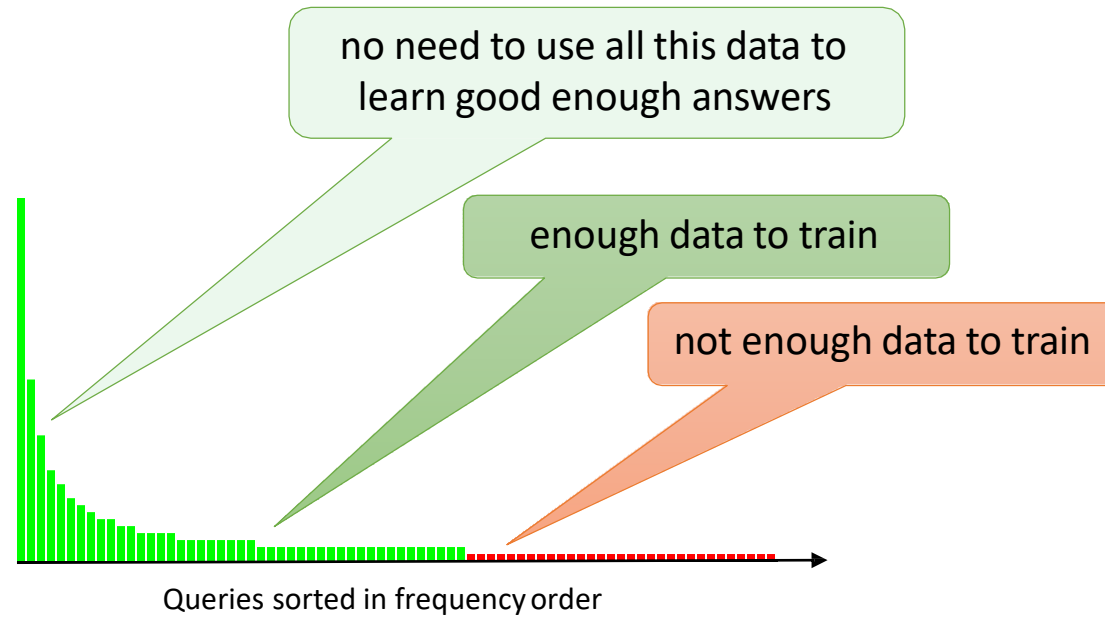How does this help ? Average accuracy is all I care about !

Training time will double, I can't afford that…

Should I optimize a different criterion?

**So what would you do ?**

# Harness Data, the scalable way



- No need to consider all examples of already known queries.
- Best is to focus on queries near the boundary of the known area.
- Curriculum learning and active learning come naturally in this context.
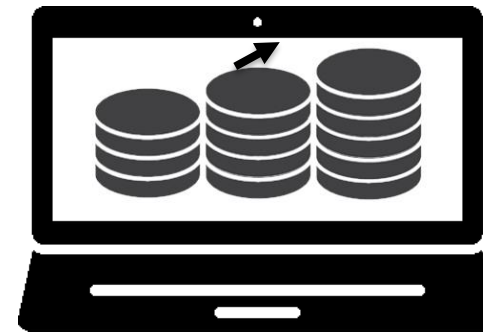- Scalability gains across the board.

# Content

Tradeoffs ✓

More Data ? ✓

Labelling

Scaling Tricks

# Labelling

# Data Augmentation

Cheap ways to get more labelled data !
Can you find some other ways ?

**Dataset**

| | |
|---|---|
| **Image** | |
| **Audio** | |
| **Text** | |
| **Video** | |

**Cheap Trick**

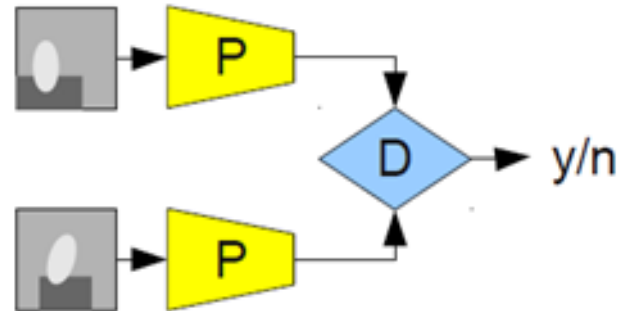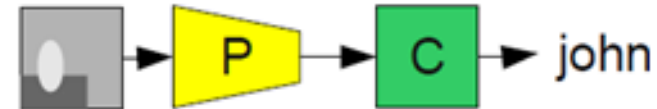| |
|---|
| Translation, rotation, change scale, add noise... |
| noise, pitch... |
| Replace word with synonym, syntax change, add/remove words, (shuffle) |
| No need to label all images ; images from same sequence should have same label. |

# Transfer Learning

**Idea**
- Learn a representation on auxilliary task with cheap labels
- Finetune for your task

Example:  face recognition

Or skip and use  a pretrained model!
- Image processing
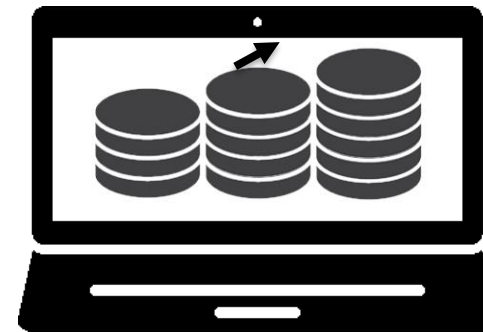- Words embeddings



(Matt Miller, NECLA, 2006)

# Content



Tradeoffs

More Data ?

Labelling

Scaling Tricks