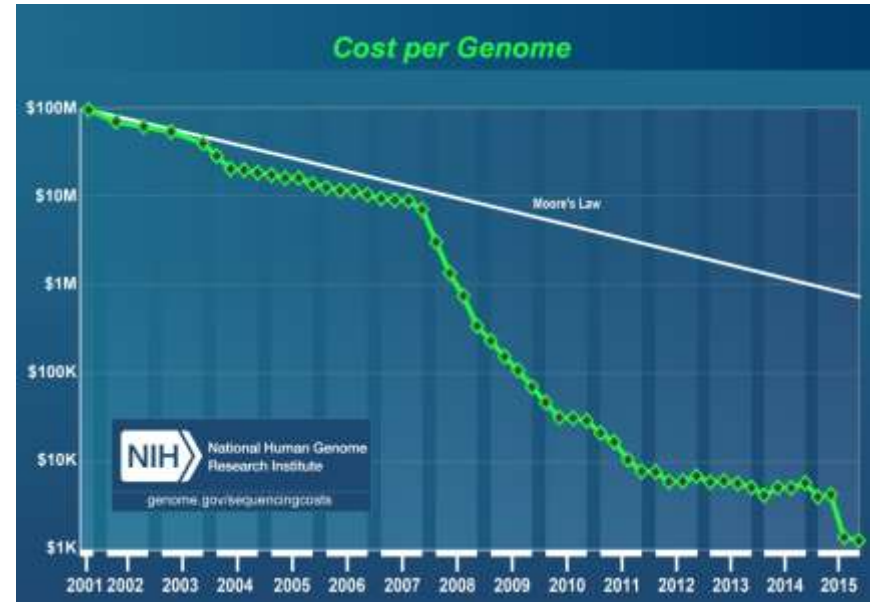


ThinkFast: Scaling Machine Learning to Modern Demands

Hristo Paskov

The Genomic Data Deluge

- Precision Medicine Initiative: sequence **1,000,000** genomes
 - **\$215 Million** in 2015
 - Pilot study
 - Outputs **10-50 GB/person**



❖ How do we analyze all of this data to **drive progress**?

Massive Data Sources

Bioinformatics



100K Genomes



News



The New York Times

Bloomberg
NEWS

Social Media



WIKIPEDIA
The Free Encyclopedia

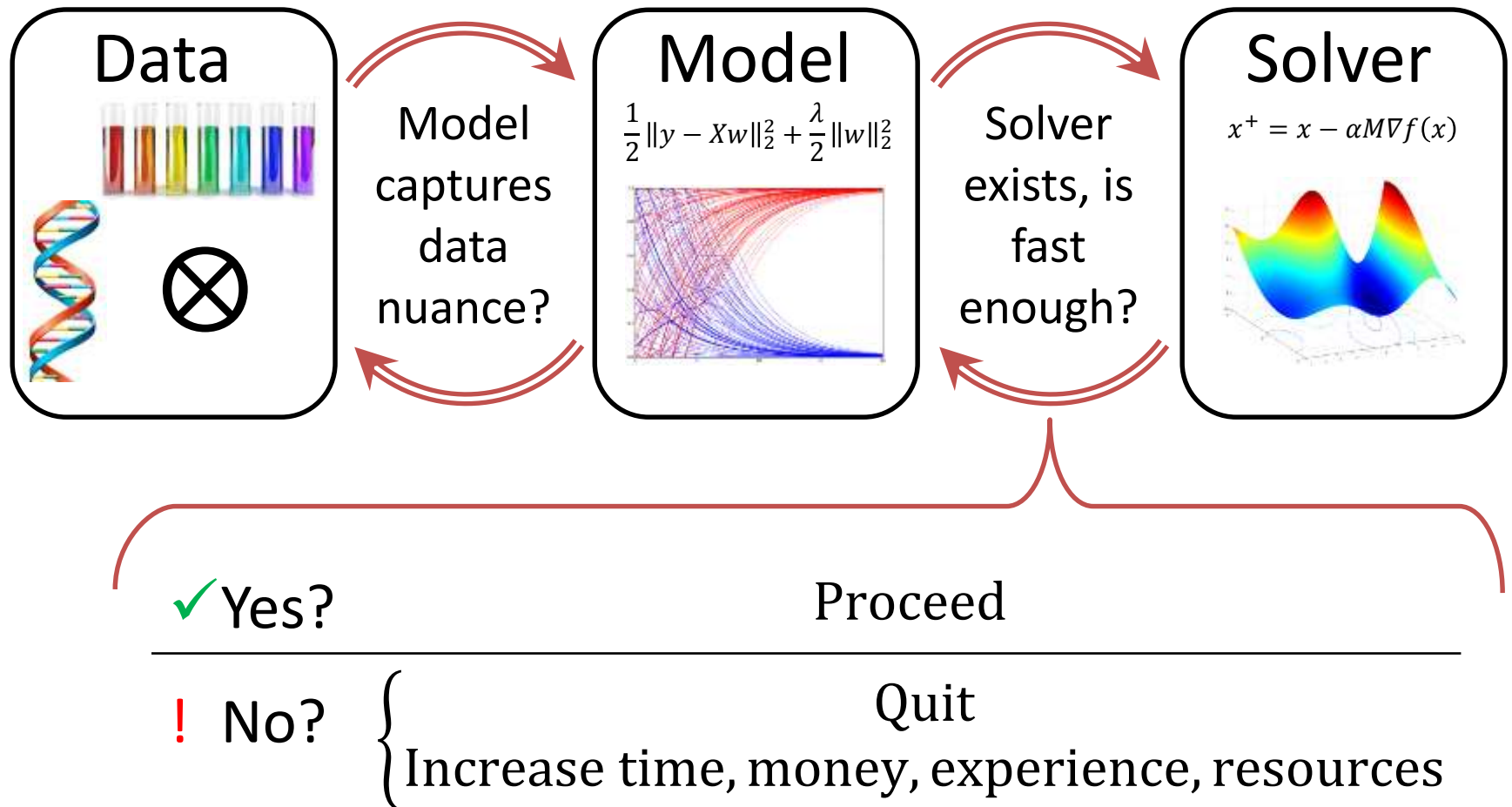
spinn3r

eCommerce

amazon

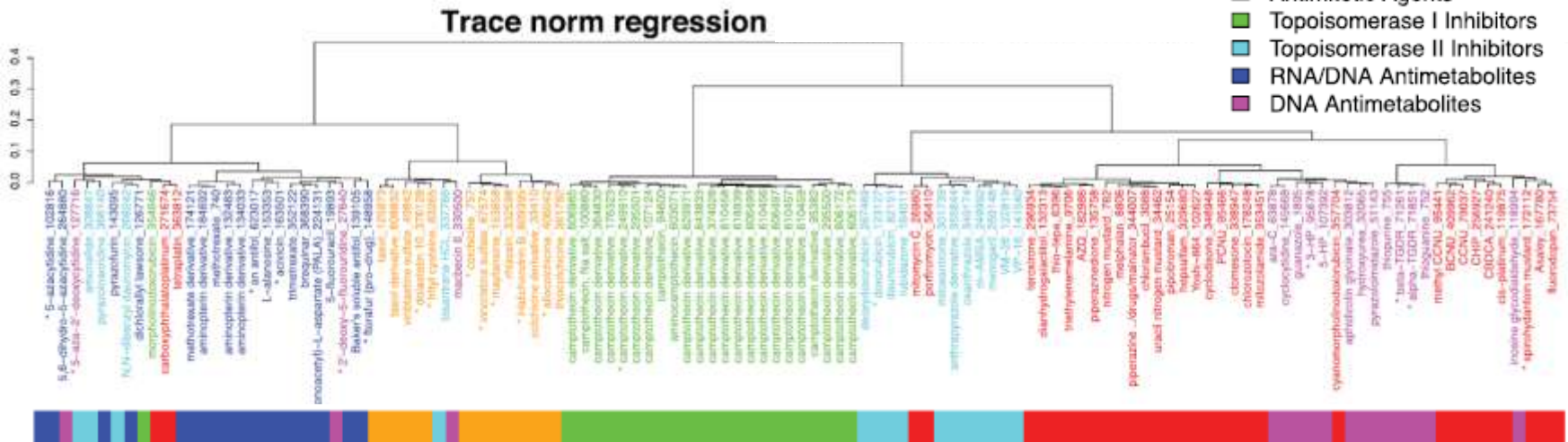
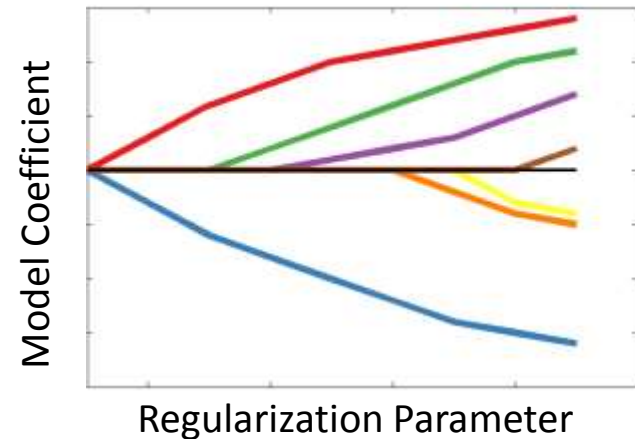


The Analysis Refinement Cycle



More Than Just Training Models

- Regularization paths
- Model risk assessment
- Interpretability

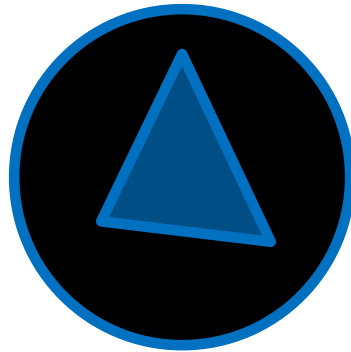


Brief History of Statistical Learning

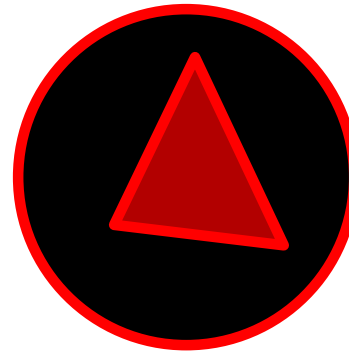
Simple
Models



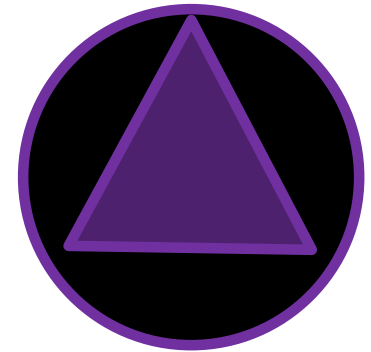
Kernel
Methods



Trees &
Ensembles



Structured
Regularization



Interpretability & Statistical Guarantees



Structured Regularization

$$\min_{\beta \in \mathbb{R}^d} L(X\beta) + \lambda R(\beta)$$

Losses

Regression

Classification

Ranking

Motif Finding

Matrix Factorization

Feature Embedding

Data Imputation

...

Regularizers

Sparsity

Spatial/ Temporal /

Manifold Structure

Group Structure

Hierarchical Structure

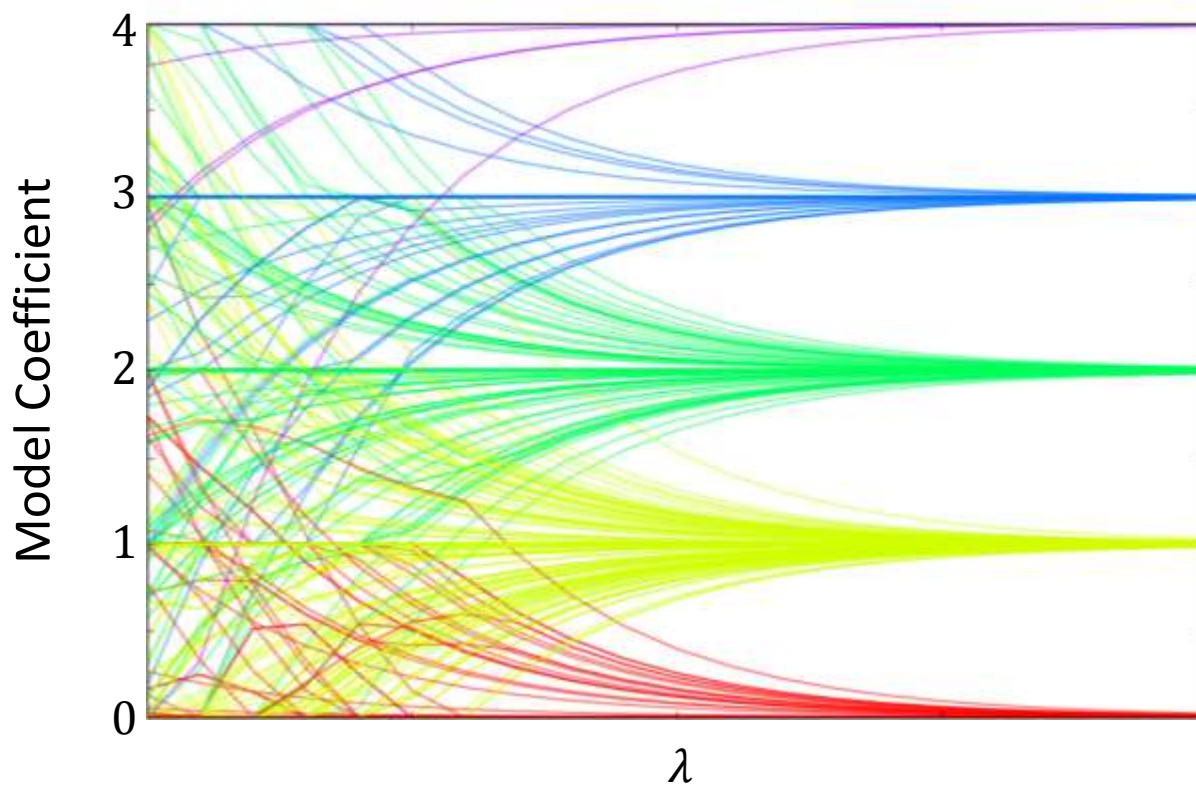
Structured & Unstructured

Multitask Learning

...

The Lasso's Combinatorial Side

$$\min_{\beta \in \mathbb{R}^d} L(y - X\beta) + \lambda \|\beta\|_1$$



The Database Perspective

$$\min_{\beta \in \mathbb{R}^d} L(y - X\beta) + \lambda \|\beta\|_1$$

$$-X^T \partial_{y-X\beta} L(y - X\beta) + \lambda \partial_{\beta} \|\beta\|_1$$

The Database Perspective

$$-\mathbf{X}^T \partial_{\mathbf{y} - \mathbf{X}\boldsymbol{\beta}} L(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \partial_{\boldsymbol{\beta}} \|\boldsymbol{\beta}\|_1$$

Feature & label storage



The Database Perspective

$$-X^T \partial_{y-X\beta} L(y - X\beta) + \lambda \partial_{\beta} \|\beta\|_1$$

Data access operations

$$u = y - X\beta$$

$$v = \partial_u L(u)$$

$$w = X^T v$$

Feature & label storage



The Database Perspective

$$-X^T \partial_{y-X\beta} L(y - X\beta) + \lambda \partial_{\beta} \|\beta\|_1$$

ML “Query Language”

$$\min_{\beta \in \mathbb{R}^d} L(y - X\beta) + \lambda \|\beta\|_1$$

Data access operations

$$\begin{aligned} u &= y - X\beta \\ v &= \partial_u L(u) \\ w &= X^T v \end{aligned}$$

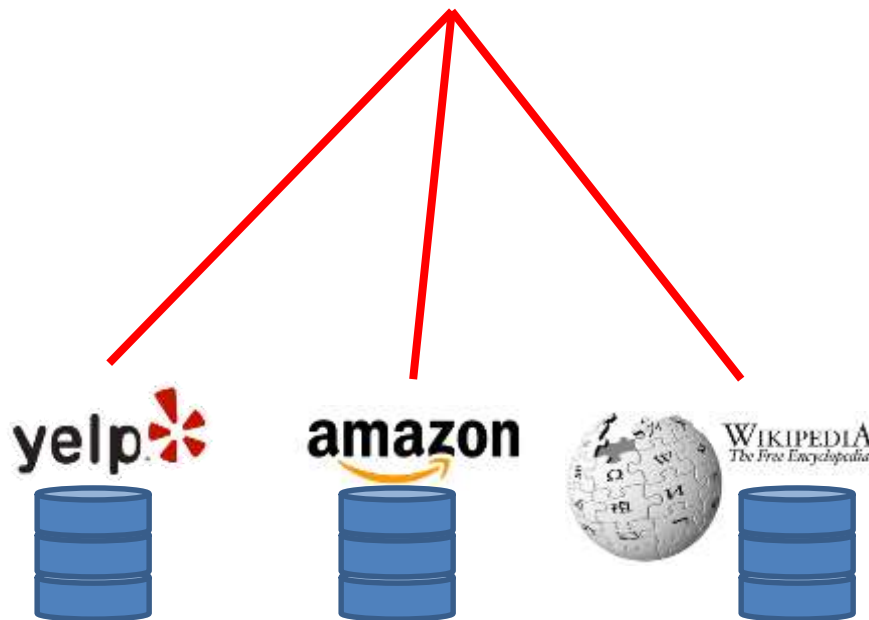
Feature & label storage



The Database Perspective

$$\min_{\beta_1, \beta_2, \beta_3 \in \mathbb{R}^d} \sum_{t=1}^3 [L_t(y_t - X_t \beta_t) + \lambda_t R_t(\beta_t)]$$

$$+ \omega \| [\beta_1 \quad \beta_2 \quad \beta_3] \|_*$$



The Database Perspective

ML “Query Language”

Data access operations

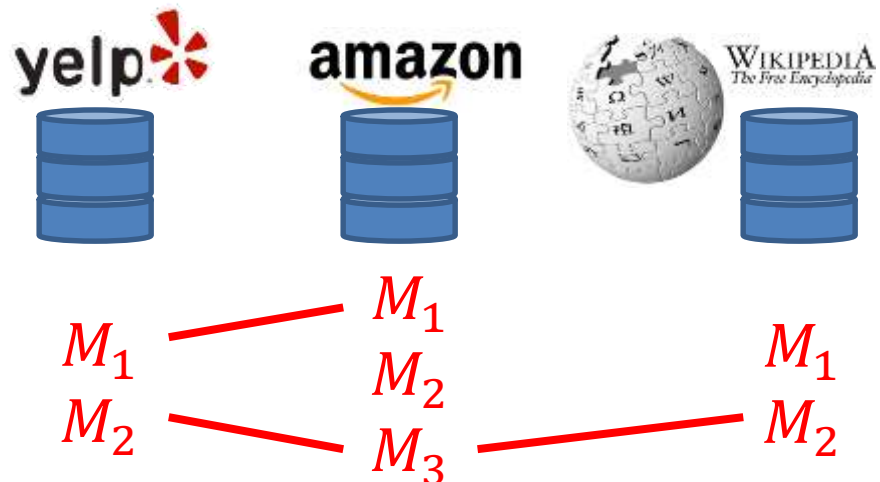
Feature, label and
model storage

$$\min_{\beta \in \mathbb{R}^d} L(y - X\beta) + \lambda \|\beta\|_1$$

$$u = y - X\beta$$

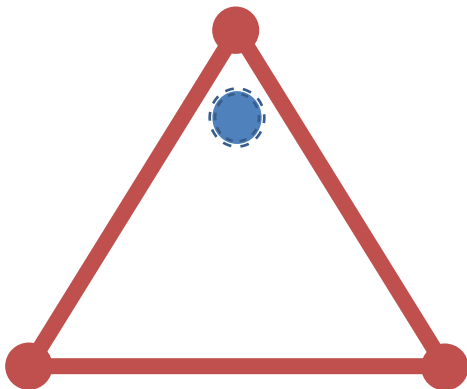
$$v = \partial_u L(u)$$

$$w = X^T v$$



The Database Perspective

Mathematical
Structure



Processing

Memory

$$\min_{\beta \in \mathbb{R}^d} L(y - X\beta) + \lambda \|\beta\|_1$$

$$u = y - X\beta$$

$$v = \partial_u L(u)$$

$$w = X^T v$$

yelp



amazon



WIKIPEDIA
The Free Encyclopedia



M_1

M_2

M_1

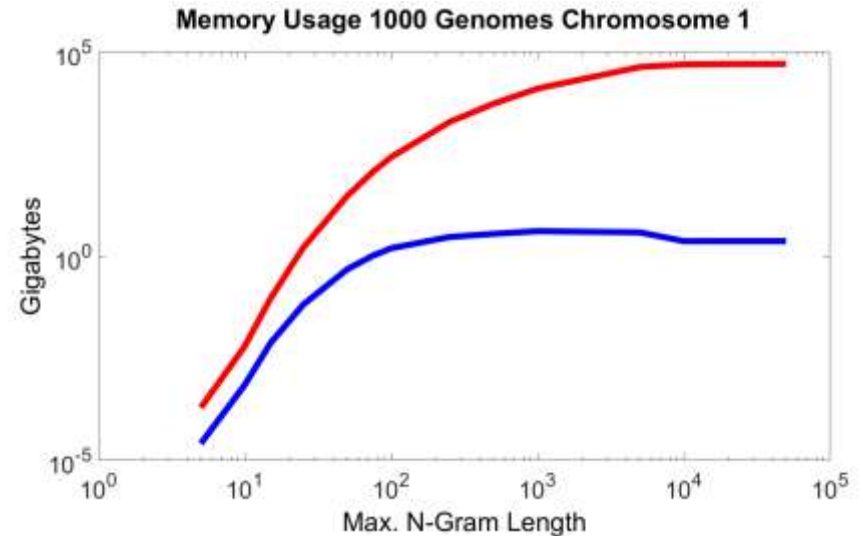
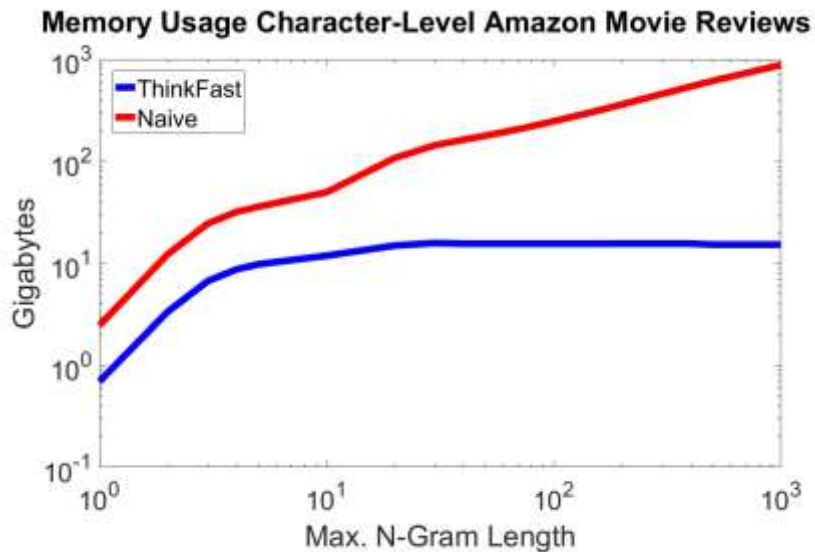
M_2

M_3

M_1

M_2

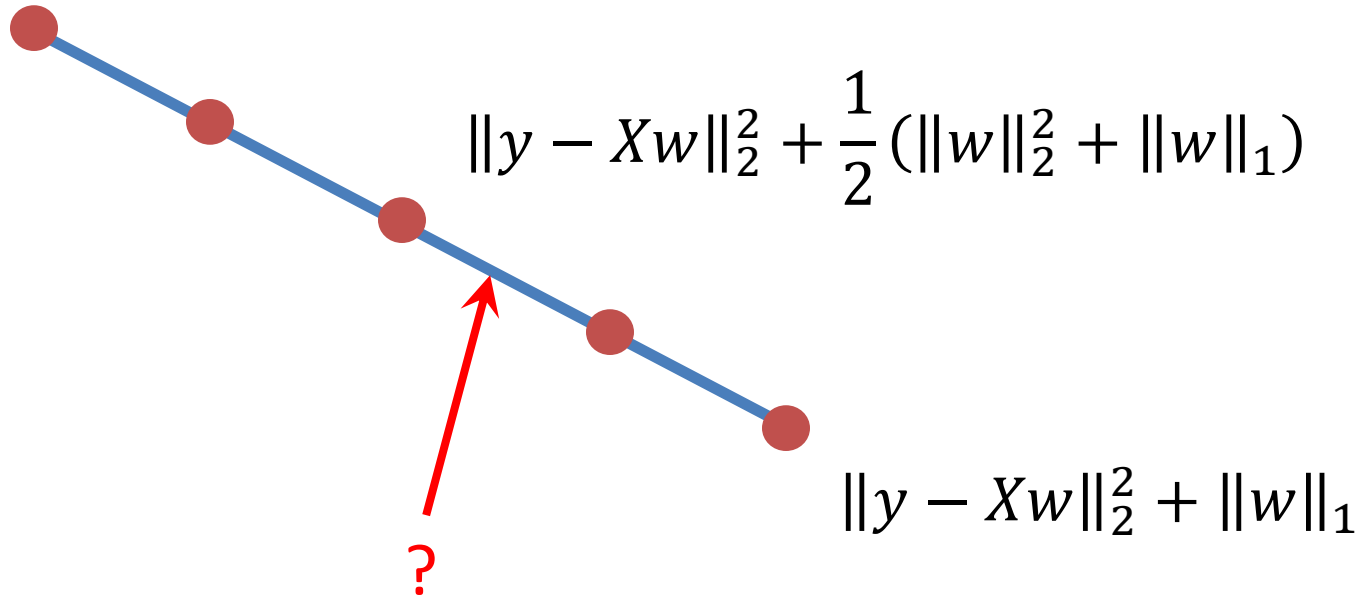
Efficient Feature Storage



“Query Language” Optimization

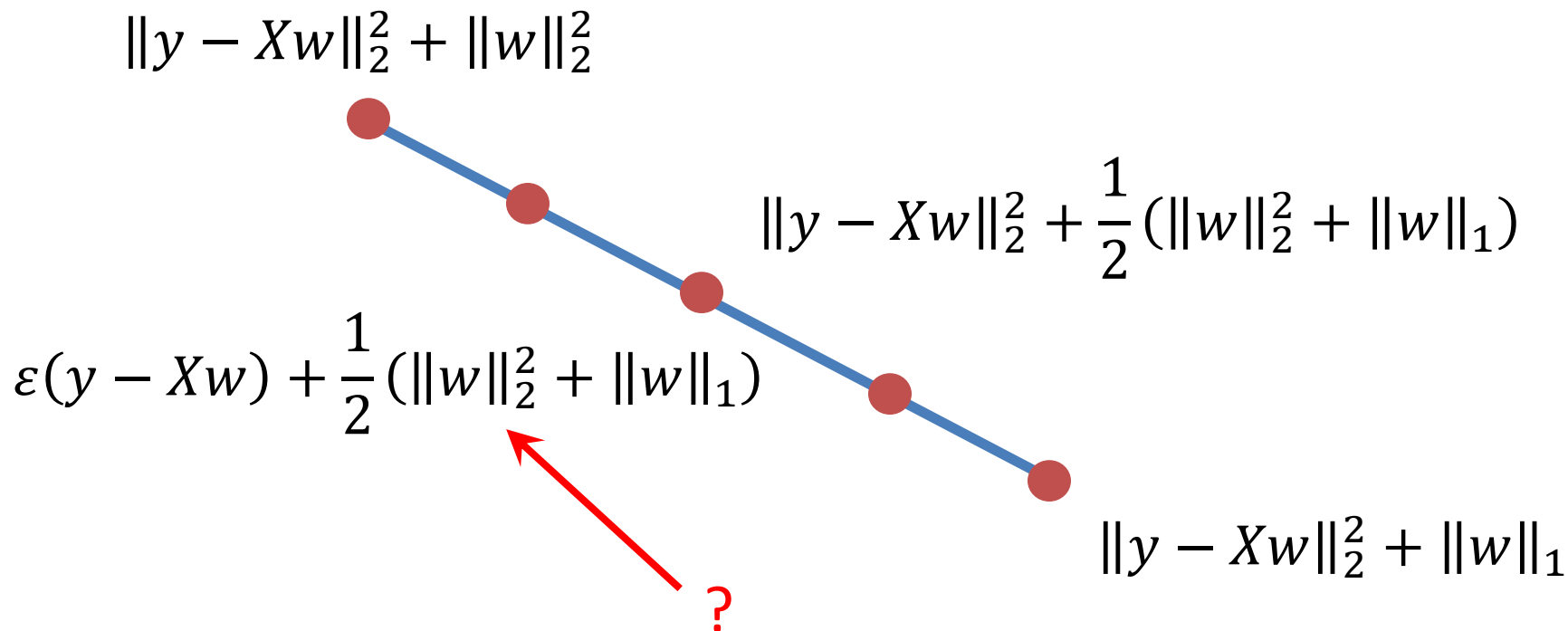
- Static analysis

$$\|y - Xw\|_2^2 + \|w\|_2^2$$

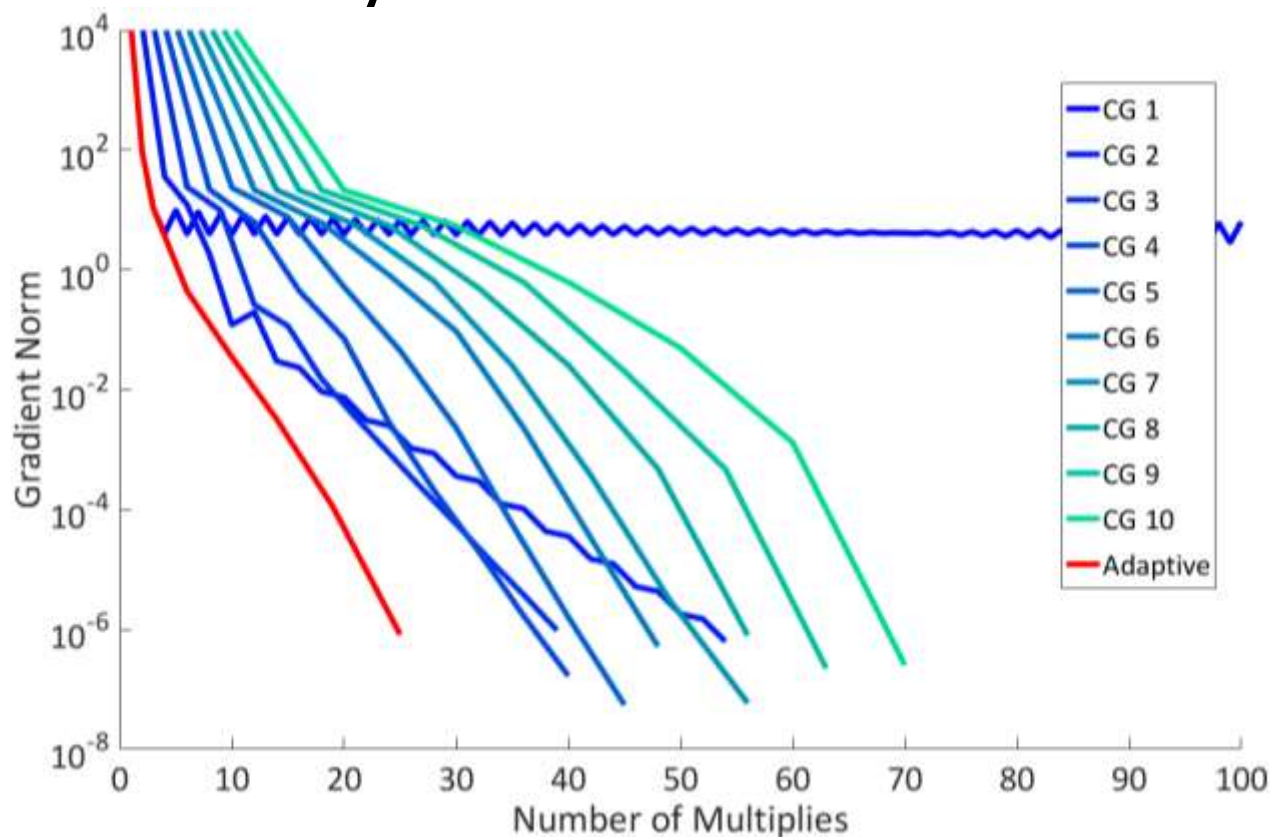


“Query Language” Optimization

- Static analysis



- Static analysis
- Runtime analysis



Some Bioinformatics Applications

- Personalized medicine, Memorial Sloan Kettering Cancer Center
 - 35% accuracy improvement over state-of-the-art
- Metagenomic binning and DNA quality assessment, Stanford School of Medicine
 - Previously unsolved problem
- Toxicogenomic analysis, Stanford University
 - Improved on state-of-the-art results

Upcoming

- Massive scale character level sentiment and text analysis on Amazon data
 - Billions of features, hours to solve a model
 - Efficient multitask learning
- Characterize the global limitations of learning word structure
 - Devise provably more efficient regularizers for uncovering structure