

# One-Pass Ranking Models for Low-Latency Product Recommendations

Martin Saveski  
@msaveski

MIT  
(Amazon Berlin)

# One-Pass Ranking Models for Low-Latency Product Recommendations

Amazon Machine Learning Team, Berlin



Antonino Freno



Rodolphe Jenatton

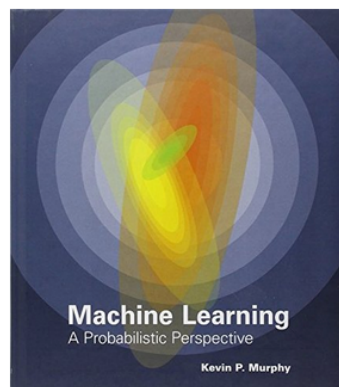


Cédric Archambeau

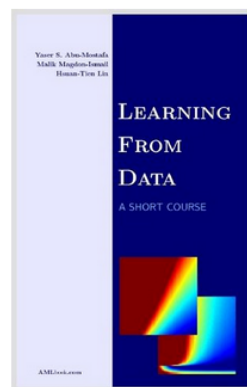
# Product Recommendations

## Customers Who Bought This Item Also Bought

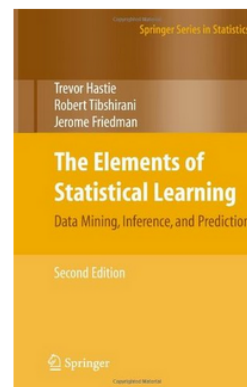
Page 1 of 20



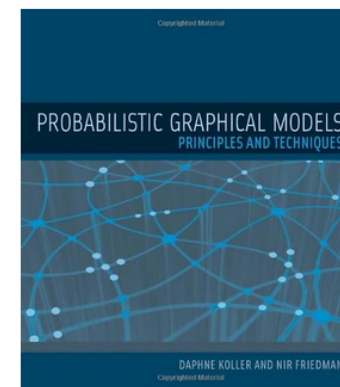
**Machine Learning: A Probabilistic Perspective**  
(Adaptive Computation and  
› Kevin P. Murphy  
★★★★☆ 46  
Hardcover  
\$76.97 ✓Prime



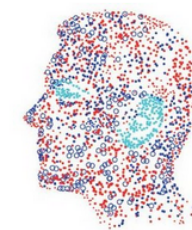
**Learning From Data**  
› Yaser S. Abu-Mostafa  
★★★★☆ 88  
#1 Best Seller in Computer  
Neural Networks  
Hardcover



**The Elements of Statistical Learning: Data Mining, Inference, and Prediction,**  
Trevor Hastie  
★★★★☆ 49  
Hardcover  
\$70.40 ✓Prime



**Probabilistic Graphical Models: Principles and Techniques (Adaptive**  
› Daphne Koller  
★★★★☆ 28  
Hardcover  
\$97.03 ✓Prime



**Machine Learning**  
The Art and Science of Algorithms that Make Sense of Data  
Peter Flach  
★★★★☆ 17  
Paperback  
\$51.60 ✓Prime



# Product Recommendations

Constraints

# Product Recommendations

## Constraints

1. Large # of examples  
Large # of features

# Product Recommendations

## Constraints

1. Large # of examples  
Large # of features
2. Drifting distribution

# Product Recommendations

## Constraints

1. Large # of examples  
Large # of features
2. Drifting distribution
3. Real-time ranking  
( $< \text{few ms}$ )

# Product Recommendations

## Constraints

1. Large # of examples  
Large # of features → Small memory footprint
2. Drifting distribution
3. Real-time ranking  
(<few ms)



# Product Recommendations

## Constraints

1. Large # of examples  
Large # of features → Small memory footprint
2. Drifting distribution → Fast training time
3. Real-time ranking  
(<few ms)

# Product Recommendations

## Constraints

1. Large # of examples  
Large # of features → Small memory footprint
2. Drifting distribution → Fast training time
3. Real-time ranking  
(<few ms) → Low prediction latency

# Our approach

## Product Recommendations

Small memory footprint

Fast training time

Low prediction latency

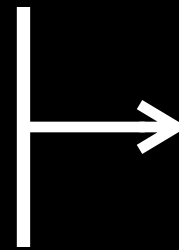
# Our approach

## Product Recommendations

Small memory footprint

Fast training time

Low prediction latency



Stochastic optimization

One pass learning

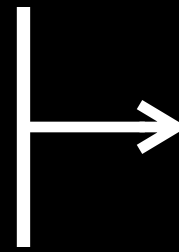
# Our approach

## Product Recommendations

Small memory footprint

Fast training time

Low prediction latency



Stochastic optimization

One pass learning

Sparse models

# Learning Ranking Functions

# Learning Ranking Functions

Three broad families of models

1. Pointwise (Logistic regression)
2. Pairwise (RankSVM)
3. Listwise (ListNet)

# Learning Ranking Functions

Three broad families of models

1. Pointwise (Logistic regression)
2. Pairwise (RankSVM)
3. Listwise (ListNet)

Loss functions

- Evaluation functions (NDCG)
- Surrogate functions



# Loss Function

Lambda Rank (Burges et al., 2007)

# Loss Function

Lambda Rank (Burges et al., 2007)

	Product 1	Product 2	Product 3	Product 4
<b>X</b> : Features	<b>x<sub>1</sub></b>	<b>x<sub>2</sub></b>	<b>x<sub>3</sub></b>	<b>x<sub>4</sub></b>
<b>r</b> : Ground-truth Rank	1	1	2	3

# Loss Function

Lambda Rank (Burges et al., 2007)

	Product 1	Product 2	Product 3	Product 4
$\mathbf{X}$ : Features	$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{x}_3$	$\mathbf{x}_4$
$\mathbf{r}$ : Ground-truth Rank	1	1	2	3
		$i$		$j$

# Loss Function

Lambda Rank (Burges et al., 2007)

	Product 1	Product 2	Product 3	Product 4
$\mathbf{X}$ : Features	$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{x}_3$	$\mathbf{x}_4$
$\mathbf{r}$ : Ground-truth Rank	1	1	2	3
		$i$		$j$

**Importance of sorting  $i$  and  $j$  correctly**

$$\Delta\mathcal{M} = \mathcal{M}(\mathbf{r}) - \mathcal{M}(\mathbf{r}_{i/j})$$

# Loss Function

Lambda Rank (Burges et al., 2007)

	Product 1	Product 2	Product 3	Product 4
$\mathbf{X}$ : Features	$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{x}_3$	$\mathbf{x}_4$
$\mathbf{r}$ : Ground-truth Rank	1	1	2	3
		$i$		$j$

**Importance of sorting  $i$  and  $j$  correctly**

$$\Delta \mathcal{M} = \mathcal{M}(\mathbf{r}) - \mathcal{M}(\mathbf{r}_{i/j})$$

**Difference in scores**

$$\Delta S = \max\{0, \mathbf{w}^T \mathbf{x}_j - \mathbf{w}^T \mathbf{x}_i\}$$

# Loss Function

Lambda Rank (Burges et al., 2007)

	Product 1	Product 2	Product 3	Product 4
$\mathbf{X}$ : Features	$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{x}_3$	$\mathbf{x}_4$
$\mathbf{r}$ : Ground-truth Rank	1	1	2	3
		$i$		$j$

**Importance of sorting  $i$  and  $j$  correctly**

$$\Delta\mathcal{M} = \mathcal{M}(\mathbf{r}) - \mathcal{M}(\mathbf{r}_{i/j})$$

**Difference in scores**

$$\Delta S = \max\{0, \mathbf{w}^T \mathbf{x}_j - \mathbf{w}^T \mathbf{x}_i\}$$

**Loss**

$$L(\mathbf{X}; \mathbf{w}) = \sum_{\mathbf{r}_i \leq \mathbf{r}_j} \Delta\mathcal{M} \cdot \Delta S$$

# ElasticRank

Introducing Sparsity

Adding  $l_1$  and  $l_2$  penalties

$$L^*(\mathbf{X}, \mathbf{w}) = L(\mathbf{X}, \mathbf{w}) + \lambda_1 ||\mathbf{w}||_1 + \frac{1}{2} \lambda_2 ||\mathbf{w}||_2^2$$

# ElasticRank

Introducing Sparsity

Adding  $l_1$  and  $l_2$  penalties

$$L^*(\mathbf{X}, \mathbf{w}) = L(\mathbf{X}, \mathbf{w}) + \lambda_1 ||\mathbf{w}||_1 + \frac{1}{2} \lambda_2 ||\mathbf{w}||_2^2$$

Both  $\lambda_1$  and  $\lambda_2$  control model complexity



# ElasticRank

Introducing Sparsity

Adding  $l_1$  and  $l_2$  penalties

$$L^*(\mathbf{X}, \mathbf{w}) = L(\mathbf{X}, \mathbf{w}) + \lambda_1 ||\mathbf{w}||_1 + \frac{1}{2} \lambda_2 ||\mathbf{w}||_2^2$$

Both  $\lambda_1$  and  $\lambda_2$  control model complexity

- $\lambda_1$  trades-off sparsity and performance

# ElasticRank

Introducing Sparsity

Adding  $l_1$  and  $l_2$  penalties

$$L^*(\mathbf{X}, \mathbf{w}) = L(\mathbf{X}, \mathbf{w}) + \lambda_1 ||\mathbf{w}||_1 + \frac{1}{2} \lambda_2 ||\mathbf{w}||_2^2$$

Both  $\lambda_1$  and  $\lambda_2$  control model complexity

- $\lambda_1$  trades-off sparsity and performance
- $\lambda_2$  adds strong convexity & improves convergence

# Optimization Algorithms

Extensions of Stochastic Gradient Descent

# Optimization Algorithms

Extensions of Stochastic Gradient Descent

**FOBOS** Forward-Backward Splitting (Duchi, 2009)

1. Gradient step
2. Proximal step involving the regularization

# Optimization Algorithms

Extensions of Stochastic Gradient Descent

**FOBOS** Forward-Backward Splitting (Duchi, 2009)

1. Gradient step
2. Proximal step involving the regularization

**RDA** Regularized Dual Averaging (Xiao, 2010)

- Keeps a running average of all past gradients
- Solves a proximal step using the average

# Optimization Algorithms

Extensions of Stochastic Gradient Descent

**FOBOS** Forward-Backward Splitting (Duchi, 2009)

1. Gradient step
2. Proximal step involving the regularization

**RDA** Regularized Dual Averaging (Xiao, 2010)

- Keeps a running average of all past gradients
- Solves a proximal step using the average

**pSGD** Pruned Stochastic Gradient Descent

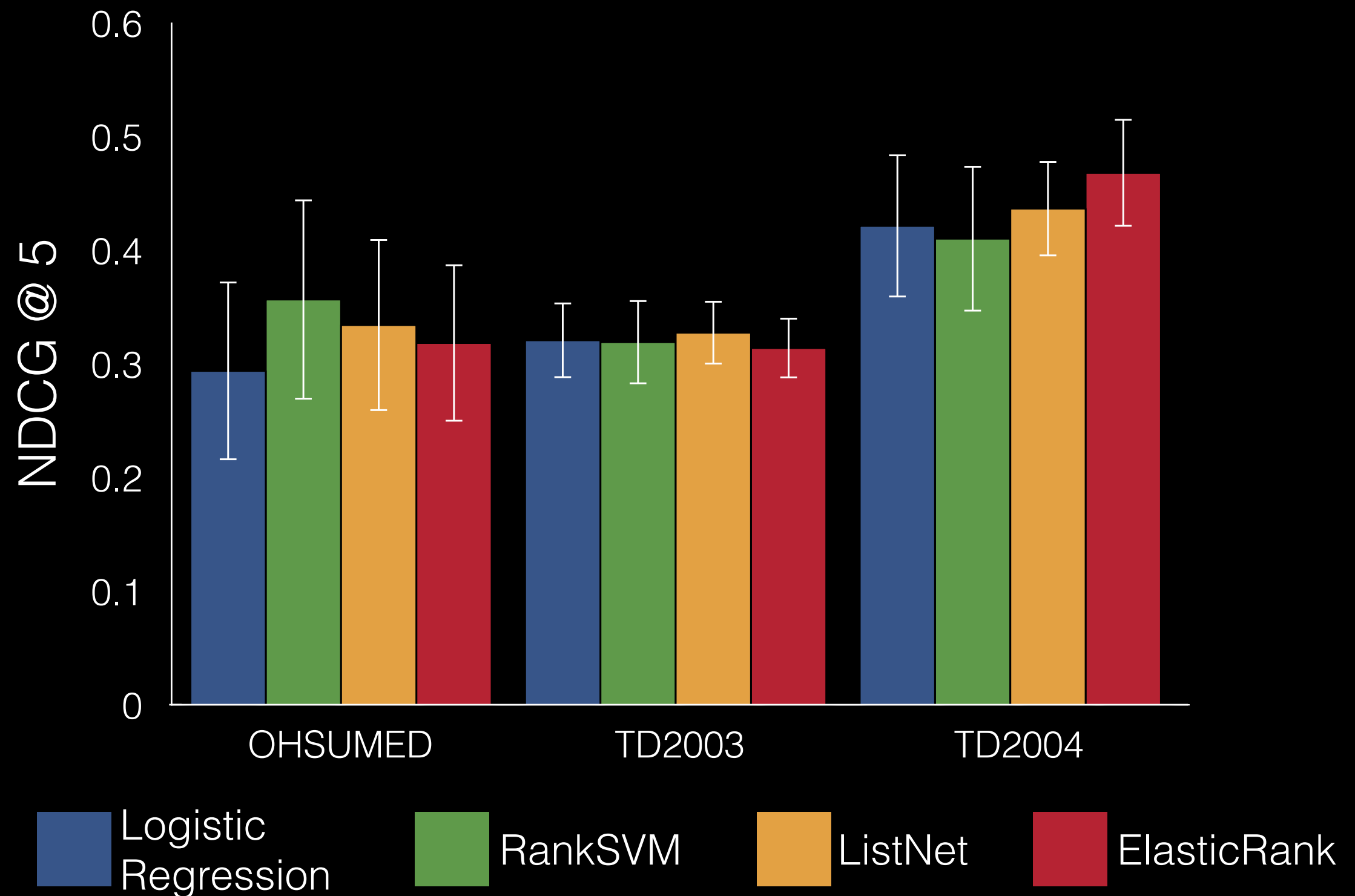
- Prunes every  $k$  gradient steps
- If  $|w_i| < \theta \Rightarrow w_i = 0$

# Hyper-parameter Optimization

- **Turn-key** inference
- Automatic adjustment of hyper-parameters
- Bayesian Approach (Snoek, Larochelle, Adams; 2012)
  - Gaussian Process
  - Thomson Sampling

# LETOR Experiments

ElasticRank is comparable with state-of-the-art models





# Amazon.com Experiments

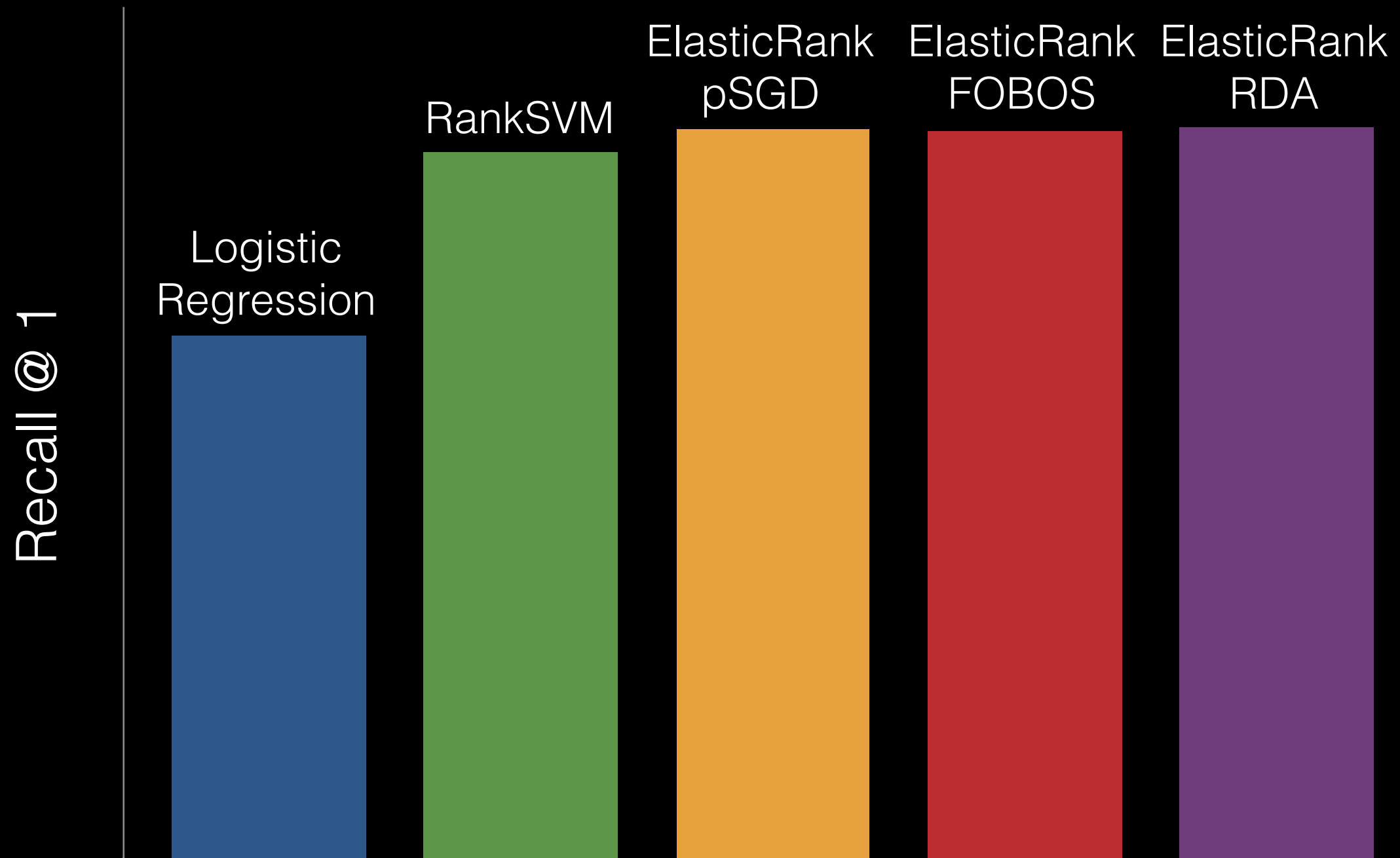
## Experimental Setup

- # examples  $\approx$  millions
- # features  $\approx$  thousands (millions of dimensions)
- Purchase logs from contiguous time interval



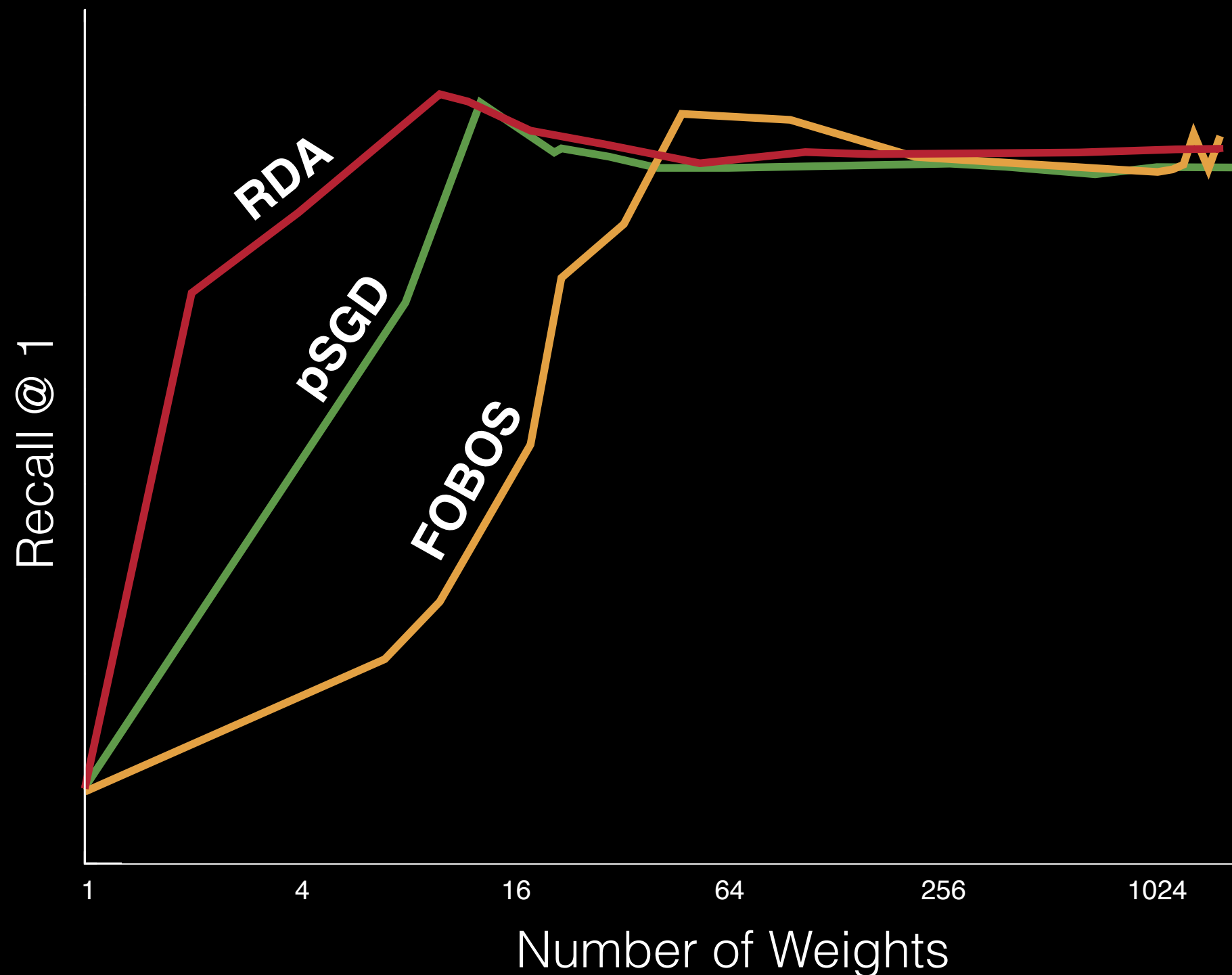
# Experimental Results

ElasticRank performs best

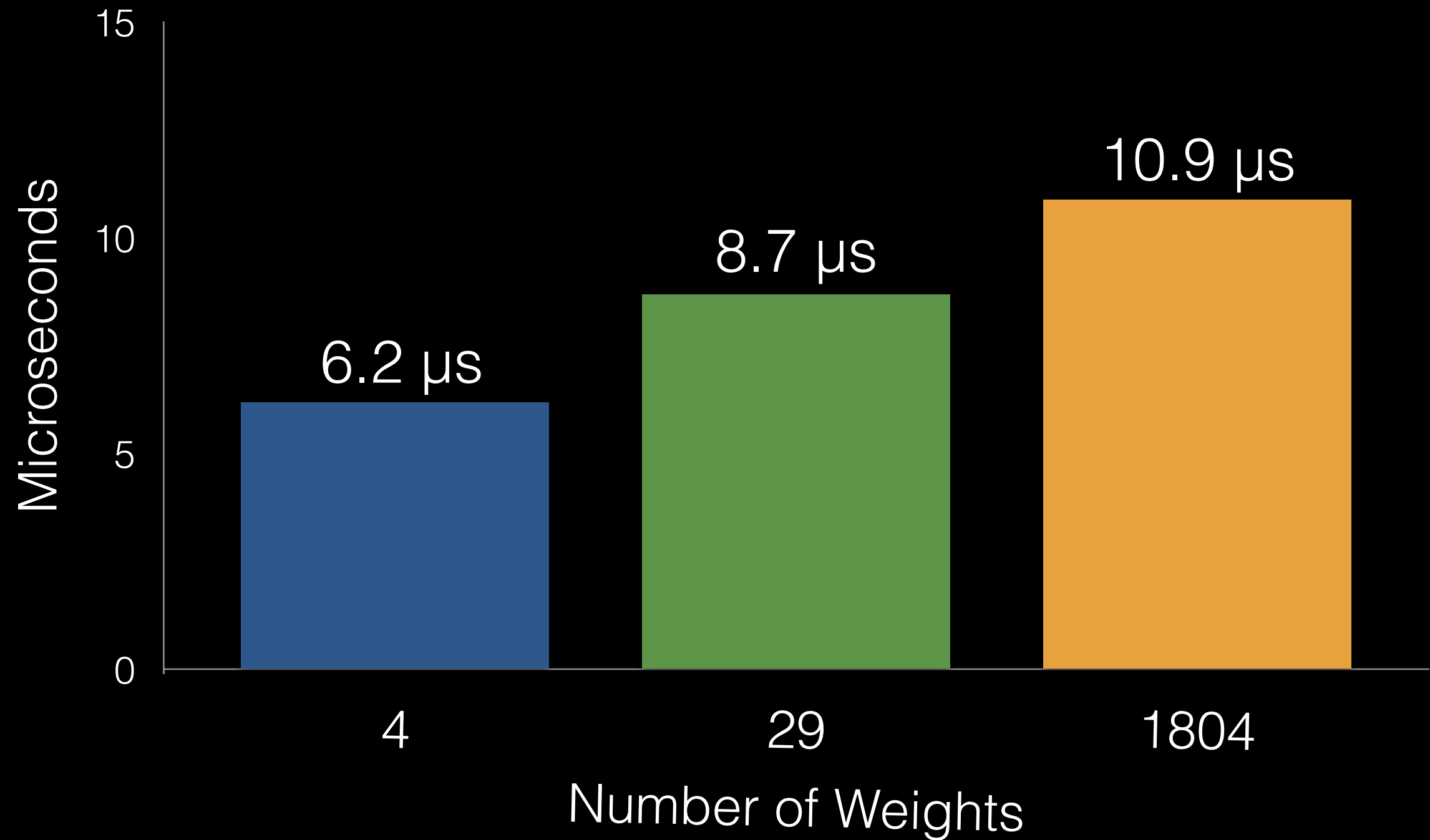


# Sparsity vs Performance

RDA achieves the best trade-off



# Prediction Time



# Contributions

How to learn ranking functions with

- Single pass
- Small memory footprint
- Sparse

**WITHOUT** sacrificing performance

# References

- C. J. C. Burges, R. Ragno, and Q. V. Le. *Learning to rank with nonsmooth cost functions*. In Advances in Neural Information Processing Systems (NIPS), 2006.
- J. C. Duchi and Y. Singer. *Efficient online and batch learning using forward backward splitting*. Journal of Machine Learning Research (JMLR), 2009.
- L. Xiao. *Dual Averaging Methods for Regularized Stochastic Learning and Online Optimization*. Journal of Machine Learning Research (JMLR), 2010.
- J. Snoek, H. Larochelle, and R. P. Adams. *Practical bayesian optimization of machine learning algorithms*. In Advances in Neural Information Processing Systems (NIPS), 2012.

# One-Pass Ranking Models for Low-Latency Product Recommendations

Martin Saveski  
@msaveski

MIT  
(Amazon Berlin)