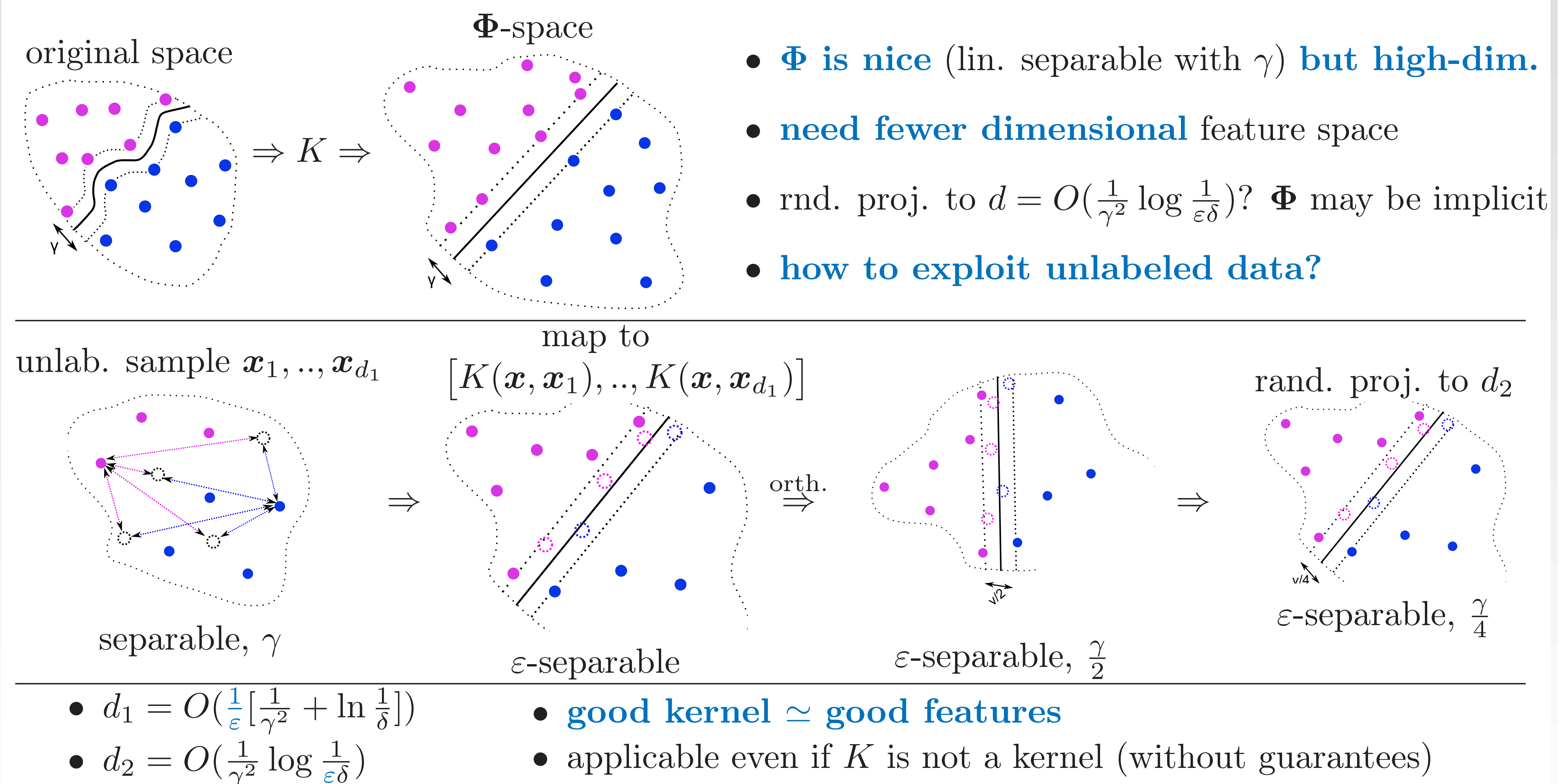


HIGHLIGHTS

WWW – millions of features

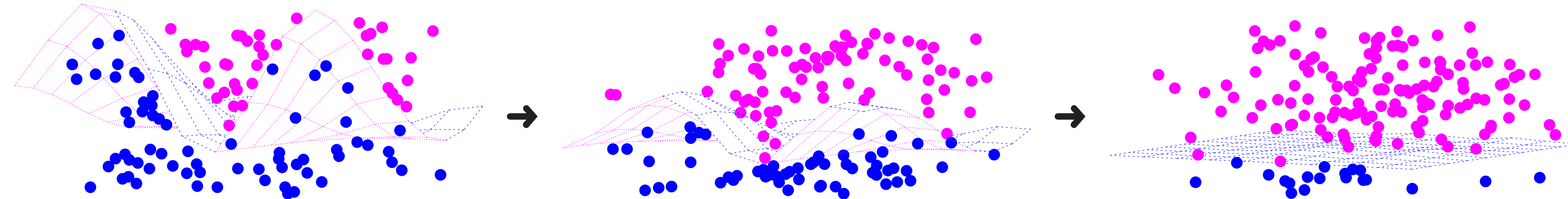
- **too much data:**
 - impossible to keep all on disk
 - necessary to learn on it
 - learning task unknown beforehand
- reduce data to few **informative** features
- **reduced data must permit learning**
- ★ **2nd and 3rd place** in Semi-Supervised Feature Learning Challenge (SSFL)
- ★ **faithful submissions:**
 - using unlabeled data
 - not using test data for self-tuning
 - no “single feature” trick (for 3rd place)

KERNELS AS FEATURES (USED HERE AS A “BLACK-BOX” PROCEDURE) [BALCAN ET AL., 2004]

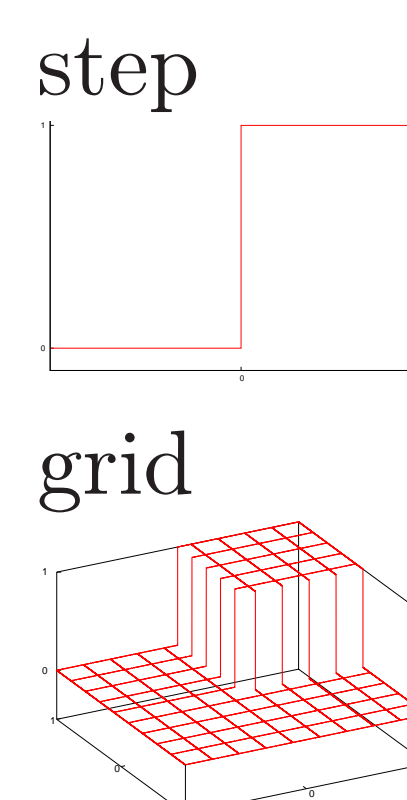


CONTRIB: TUNE KERNELS WITH BOOSTING & NEURAL NETWORKS

I – Non-linear feature **transform with RankBoost** to make data separable by a hyperplane



- Seeking scoring function as: $H(x) = \sum_{t=1}^T \alpha_t h_t(x)$
- Optimizing pair-wise loss, equivalent to AUC: $L = \sum_{y_i < y_j} [H(x_i) > H(x_j)]$
- Weak rankers $h(x)$: decision stumps and grids (trees of depth 2)
- Use **weak outputs as new features** $\Phi(x) = (\alpha_1 h_1(x), \dots, \alpha_T h_T(x))$
- H can be **viewed as linear discrim. rule** $H(x) = \langle w, \Phi(x) \rangle$ for $w = \bar{1}$
- T may be too big to store training $\Phi(x_i)$ + **want to exploit unlabeled data** \Rightarrow
- **define kernel** as $K(x_1, x_2) = \langle \Phi(x_1), \Phi(x_2) \rangle \Rightarrow$ **use “black-box”**

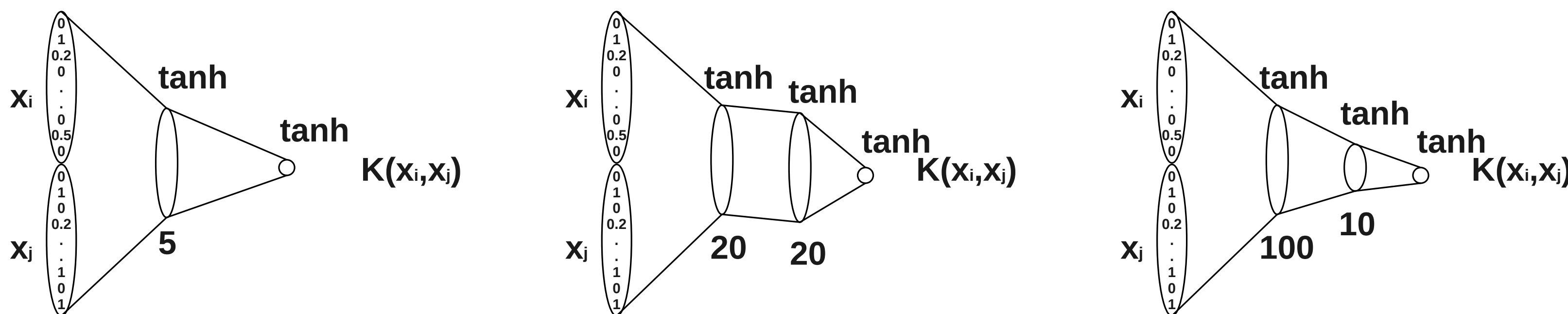


II – Neural network K to optimize **kernel alignment** with perfect kernel $K'(x_i, x_j) = y_i y_j$

$$A(K, K') = K \cdot K' / \sqrt{(K \cdot K)(K' \cdot K')} \quad K \cdot K' = \sum_{i,j} K(x_i, x_j) y_i y_j$$

- Optimize **quadratic error**: $\sum_{x_i, x_j} (K(x_i, x_j) - y_i y_j)^2$ with stochastic gradient descend
- Normalization factor in A is **ignored**
- Φ -representation not accessible \Rightarrow **use “black-box”**

Tested configurations



STAGES & VARIANCES

A RankBoost features B projection onto unlabeled data C random projection D classifier learning

kernel	steps	I_{sample}	A→D	A→B→D	A→C→D	A→B→C→D
stump	2270	1000	0.99539	0.99280	0.99517	0.99232
stump	2270	5000	0.99539	0.99296	0.99517	0.99166
grid	2000	1000	0.99076	0.99518	0.99416	0.99378
grid	2000	5000	0.99076	0.99517	0.99416	0.99546
kernel	steps	I_{sample}	A→B→D		A→B→C→D	
stump	2270	100	0.9916±0.0015		0.9935±0.0017	
stump	2270	1000	0.9926±0.0003		0.9925±0.0004	
stump	2270	5000	0.9927±0.0001		0.9925±0.0004	
grid	2000	100	0.9950±0.0005		0.9951±0.0003	
grid	2000	1000	0.9952±0.0004		0.9950±0.0003	
grid	2000	5000	0.9950±0.0004		0.9926±0.0007	
neural 20-20	1000		0.9956±0.00003		0.9893±0.0007	
neural 100-10	1000		0.9955±0.00004		0.9886±0.0021	

Less samples \Rightarrow more variance Unexplained: Stumps: A→C→D better than A→B→D. Grids: bad A→D

SSFL CHALLENGE ON WEB DATA

- number of features – **10⁶**
 - 80% of those are binary
 - sparse (~ 115 active simultaneously)
- required output dimension – **100**
- **reduced data must permit learning**
 - fixed task: **binary linear** classif.
 - performance measure – **AUC**
- full results on poster of D. Sculley

BASELINE RESULTS

baseline	AUC
100 k-means	0.9831
1000 k-means + random projection	0.9846
1000 k-means + neural dim. red.	0.9868
RankBoost, stumps, 5000 steps	0.9961
RankBoost, stumps, 2270 steps	0.9953
RankBoost, grids, 2000 steps	0.9949
RankBoost, grids, 1150 steps	0.9949
sparse logistic regression	0.9958
sparse log. reg. + 100 k-means	0.9963
sparse log. reg. + 200 k-means	0.9963
sparse log. regression + 800 k-means	0.9962
log. regression with neural network	0.9937
log. reg. + graph smoothing	0.9949
log. reg. with NN on relabeled data	0.9847

METHOD'S RESULTS

kernel	T	I_{sample}	orth.	AUC
stump	5000	1000	no	0.9803
stump	5000	1000	yes	0.9927
stump	2270	1000	no	0.9923
stump	5000	5000	no	0.9932
stump	5000	5000	yes	0.9920
stump	2270	5000	no	0.9917
grid	2000	1000	no	0.9951
grid	2000	1000	yes	0.9938
grid	2000	5000	no	0.9955
grid	2000	5000	yes	0.9951
neural 5	1000		no	0.9895
neural 20-20	1000		no	0.9887
neural 100-10	1000		no	0.9872
neural 5	5000		no	0.9901
neural 20-20	5000		no	0.9928
neural 100-10	5000		no	0.9922