Robustness of Sketched Linear Classifiers to Adversarial Attacks

w: Linear Classifier	A: Adversary	z: Sketched Classifier
 Online linear binary classifier Learning Rate: η Convex and differentiable loss : l(·) Regularization parameter: λ 	 Aim is to induce misclassifications Targets a trained classifier at test time Budget(δ): maximum allowed perturbation δ crafted using Fast Gradient Sign 	 Properties of a sketching-based linear classifier WM-SKETCH: Reduced memory by hashing features to lower dimensions Accurate weight estimates of w

• Trained using Stochastic Gradient Descent (SGD)

Paradigm

- \mathcal{A} crafts a perturbation based on its
- SKETCH?
- count sketch R?

 \mathcal{A}

Method (FGSM)

• Adds δ to data point to move to other side



- - Sketch size k
 - Sketch depth s
 - Count sketch projection matrix R

Attack Protocol

The protocol involves choosing a surro-

FACULTY OF SCIENCE

ANANTH MAHADEVAN¹ ARPIT MERCHANT¹ YANHAO WANG² MICHAEL MATHIOUDAKIS¹ ¹ University of Helsinki, ² East China Normal University

